

1 **NONLINEAR GRADIENT MAPPINGS AND STOCHASTIC**
2 **OPTIMIZATION: A GENERAL FRAMEWORK WITH**
3 **APPLICATIONS TO HEAVY-TAIL NOISE**

4 DUŠAN JAKOVETIĆ*, DRAGANA BAJOVIĆ†, ANIT KUMAR SAHU‡, SOUMMYA KAR§,
5 NEMANJA MILOŠEVIĆ*, AND DUŠAN STAMENKOVIĆ*

6 **Abstract.** We introduce a general framework for nonlinear stochastic gradient descent (SGD)
7 for the scenarios when gradient noise exhibits heavy tails. The proposed framework subsumes several
8 popular nonlinearity choices, like clipped, normalized, signed or quantized gradient, but we also con-
9 sider novel nonlinearity choices. We establish for the considered class of methods strong convergence
10 guarantees assuming a strongly convex cost function with Lipschitz continuous gradients under very
11 general assumptions on the gradient noise. Most notably, we show that, for a nonlinearity with
12 bounded outputs and for the gradient noise that may not have finite moments of order greater than
13 one, the nonlinear SGD’s mean squared error (MSE), or equivalently, the expected cost function’s
14 optimality gap, converges to zero at rate $O(1/t^\zeta)$, $\zeta \in (0, 1)$. In contrast, for the same noise setting,
15 the linear SGD generates a sequence with unbounded variances. Furthermore, for general nonlin-
16 earities that can be decoupled component wise and a class of joint nonlinearities, we show that the
17 nonlinear SGD asymptotically (locally) achieves a $O(1/t)$ rate in the weak convergence sense and
18 explicitly quantify the corresponding asymptotic variance. Experiments show that, while our frame-
19 work is more general than existing studies of SGD under heavy-tail noise, several easy-to-implement
20 nonlinearities from our framework are competitive with state-of-the-art alternatives on real data sets
21 with heavy tail noises.

22 **Key words.** Stochastic optimization; stochastic gradient descent; nonlinear mapping; heavy-tail
23 noise; convergence rate; mean square analysis; asymptotic normality; stochastic approximation.

24 **AMS subject classifications.** 90C15, 90C25, 65K05, 62L20, 68T05

25 **1. Introduction.** Stochastic gradient descent (SGD) and its variants, e.g., [27,
26 16, 23, 35, 25, 12, 24, 7], are popular and standard methods for large scale optimization
27 and training of various machine learning models, e.g., [5, 6, 31, 8]. Recently, there have
28 been several studies that demonstrate that the gradient noise in SGD is heavy-tailed,
29 e.g., when training deep learning models [32, 17, 37].

30 Motivated by these studies, we introduce a general analytical framework for *non-*
31 *linear* SGD when the gradient evaluation is subject to a heavy-tailed noise. We combat
32 the gradient noise with a generic nonlinearity that is applied on the noisy gradient to
33 effectively reduce the noise effect. The resulting class of nonlinear methods subsumes
34 several popular choices in training machine learning models, including normalized
35 gradient descent and clipped gradient descent, e.g., [28, 36], the sign gradient, e.g.,
36 [4, 2], and (component-wise) quantized gradient, e.g., [1, 18].¹

37 We establish for the considered class of methods several results that demonstrate a
38 high degree of robustness to noise under very general assumptions on the nonlinearity

*University of Novi Sad, Faculty of Sciences, Department of Mathematics and Informatics
(dusan.jakovetic@dmi.uns.ac.rs, nmilosev@dmi.uns.ac.rs, dusan.stamenkovic@dmi.uns.ac.rs)

†University of Novi Sad, Faculty of Technical Sciences, Department of Power, Electronic and
Communication Engineering, (dbajovic@uns.ac.rs)

‡Amazon Alexa AI (anit.sahu@gmail.com)

§Department of Electrical and Computer Engineering, Carnegie Mellon University (soum-
myak@andrew.cmu.edu). The work of D. Bajovic and D. Jakovetic is partially supported by the
European Union’s Horizon 2020 Research and Innovation program under grant agreement No 957337.
The paper reflects only the view of the authors and the Commission is not responsible for any use
that may be made of the information it contains.

¹Interestingly, some of these nonlinear methods are usually introduced with a different motivation
than robustness, like, e.g., speeding up training, see, e.g., [36], or communication efficiency, [2, 4].

39 and on the gradient noise, assuming a strongly convex cost with Lipschitz continuous
 40 gradient. First, for a nonlinearity with bounded outputs (e.g., a sign, normalized,
 41 or clipped gradient) and the gradient noise that may have infinite moments of order
 42 greater than one, assuming that the noise probability density function (pdf) is sym-
 43 metric, we show that the nonlinear SGD converges almost surely to the solution, and,
 44 moreover, achieves a global $O(1/t^\zeta)$ mean squared error (MSE) convergence rate,
 45 where we explicitly quantify the degree $\zeta \in (0, 1)$. In the same setting, the linear
 46 SGD generates a sequence with unbounded variances at each iteration t . Further-
 47 more, assuming the gradient noise with finite variance, we show – for the unbounded
 48 nonlinearities that are lower bounded by a linear function – almost sure convergence
 49 and the $O(1/t)$ global MSE rate.

50 Next, for the general nonlinearities with bounded outputs that can be decoupled
 51 component-wise and a restricted class of joint nonlinearities with bounded outputs,
 52 we show under the heavy-tail noise a local (asymptotic) $O(1/t)$ rate in the weak con-
 53 vergence sense. More precisely, we show that the sequence generated by the nonlinear
 54 SGD is asymptotically normal and explicitly quantify the asymptotic variance. Fi-
 55 nally, we illustrate the results on several examples of the nonlinearity and the gradient
 56 noise pdf, highlighting and quantifying the noise regimes and the corresponding gains
 57 of the nonlinear SGD over the linear SGD scheme. In more detail, the asymptotic
 58 variance expression reveals an interesting tradeoff that the nonlinearity makes on the
 59 algorithm performance: on the one hand, the nonlinearity suppresses the noise effect
 60 to a certain degree, but on the other hand it also reduces the “useful information flow”
 61 and hence slows down convergence with respect to the noiseless case. We explicitly
 62 quantify this tradeoff and demonstrate through examples that an appropriately cho-
 63 sen nonlinearity strictly improves performance over the linear scheme in a high noise
 64 setting. Finally, we carry out numerical experiments on several real data sets that
 65 exhibit heavy tail gradient noise effects. The experiments show that, while our ana-
 66 lytical framework is more general than usual studies of SGD under heavy-tail noise,
 67 several easy-to-implement example nonlinearities of our framework – including those
 68 not previously used – are competitive with state-of-the-art alternatives.

69 Technically, for component-wise nonlinearities and the asymptotic analysis, we
 70 develop proofs based on stochastic approximation arguments, e.g., [26], following the
 71 noise and nonlinearities assumptions framework similar to [30]. The paper [30] is con-
 72 cerned with a related but different problem than ours: it considers linear estimation
 73 of a vector parameter observed through a sequence of scalar observation equations,
 74 and it is not concerned with a global MSE rate analysis that we provide here. For the
 75 MSE analysis and for the nonlinearities that cannot be expressed component-wise,
 76 like the clipped and normalized gradient, we develop novel analysis techniques.

77 There have been several works that study robustness of stochastic gradient de-
 78 scent under certain variants of heavy-tailed noises. Reference [37] consider an adap-
 79 tive gradient clipping method and establish convergence rates in expectation for the
 80 considered method under a heavy-tailed noise. For this, the authors assume that
 81 the expected value of the norm of the gradient noise raised to power α is finite, for
 82 $\alpha \in (1, 2]$. They also provide lower complexity bounds for SGD methods assuming in
 83 addition that the expected α -power of the norm of the *stochastic gradient* is finite.
 84 The paper [32] establishes convergence of the *linear* SGD assuming that the gradient
 85 noise follows a heavy-tailed α -stable distribution.

86 It is worth noting that, in addition to the MSE (expected optimality gap) results
 87 achieved here, it is also of interest to derive high probability bounds. Specifically,
 88 given a target accuracy $\epsilon > 0$ and a confidence level $1 - \beta$, $\beta \in (0, 1)$, we would like

89 to find $T(\epsilon, \beta)$ such that $f(\mathbf{x}^t) - f(\mathbf{x}^*) \leq \epsilon$ with probability at least $1 - \beta$, for all
 90 iterations $t \geq T(\epsilon, \beta)$. Application of the Markov inequality to our result $\mathbb{E}[f(\mathbf{x}^t) -$
 91 $f(\mathbf{x}^*)] = O(1/t^\zeta)$ yields, abstracting dependencies on other system parameters, a
 92 bound $T(\epsilon, \beta) \sim \frac{1}{(\beta\epsilon)^{1/\zeta}}$. This involves a strong dependence on β , on the order $1/\beta^{1/\zeta}$.
 93 Several works, e.g., [13, 14, 19, 15, 11], establish high probability bounds where $T(\epsilon, \beta)$
 94 depends *logarithmically* on β for the settings therein. For example, references [13,
 95 14] establish high probability bounds for the stochastic gradient methods therein
 96 assuming that the gradient noise has light tails (sub-Gaussian noise). The authors
 97 of [19] establish the corresponding bounds for the basic SGD and the mirror descent
 98 that utilize a gradient truncation technique. They relax the noise sub-Gaussianity
 99 assumption and assume a finite noise variance. Very recently, [15] establishes high
 100 probability bounds for accelerated SGD with a clipping nonlinearity, but assuming
 101 a finite variance of the gradient noise. Reference [11] proposes a procedure called
 102 proxBoost and establishes for the procedure high probability bounds, again assuming a
 103 finite noise variance (without the sub-Gaussianity assumption). It is highly relevant to
 104 investigate high probability bounds for the problem setting and the algorithmic class
 105 considered in this paper. Of special interest is to provide high probability bounds for
 106 a broader class of nonlinearities than the usually studied clipping-type nonlinearities;
 107 this is an interesting future work direction.

108 In summary, with respect to existing work, our framework is more general with
 109 respect to both the adopted nonlinearity in SGD and the “thickness” of the gradi-
 110 ent noise tail, assuming in addition that the noise pdf is a symmetric function. For
 111 example, current works usually assume a single choice for the nonlinearity, e.g., gradi-
 112 ent clipping, while we consider a general nonlinearity that subsumes many popular
 113 choices. Also, provided that the nonlinearity’s output is bounded (which is true for
 114 many popular choices like the clipped, signed, and normalized gradient), we establish
 115 a sublinear MSE convergence rate $O(1/t^\zeta)$ assuming only that the expected norm of
 116 the gradient noise is finite, an assumption weaker than those considered in the works
 117 of [15, 37, 11, 32]. On the other hand, we assume a strongly convex smooth cost func-
 118 tion, which is equivalent to or stronger than the assumptions made in these works.
 119 See also Examples 3.2 and 3.3. ahead for further rate comparisons with existing work.

120 The idea of employing a nonlinearity into a “baseline” linear scheme has also been
 121 used in other contexts. Most notably, several works consider nonlinear versions of the
 122 standard consensus algorithm to evaluate average of scalar values in a distributed
 123 fashion, e.g., [22, 33, 10]. The paper [22] introduces a trigonometric nonlinearity
 124 into a standard linear consensus dynamics and shows an improved dependence of the
 125 method on initial conditions. References [33] and [10] employ a general nonlinearity
 126 in the linear consensus dynamics and show that it improves the method’s resilience
 127 to additive communication noise. The authors of [34] modify the linear consensus by
 128 taking out from the averaging operation the maximal and minimal estimates among
 129 the estimates from all neighbors of a node. The above works are different from
 130 ours as they focus on the specific consensus problem that can be translated into
 131 minimizing a convex quadratic cost function in a distributed way over a generic,
 132 connected network. In contrast, we consider general strongly convex costs, and we
 133 are not directly concerned with distributed systems.

134 **Paper organization.** Section 2 describes the problem model and the nonlinear
 135 SGD framework that we assume. Section 3 and Section 4 explain our results on
 136 nonlinear SGD for component-wise and joint nonlinearities, respectively. Section 5
 137 and Section 6 then provide proofs of the corresponding results. Section 7 illustrates

138 the performance of several example methods from our nonlinear SGD framework on
 139 real data sets that have heavy-tail gradient noise. Finally, [Section 8](#) concludes the
 140 paper. Some auxiliary results and proofs are delegated to the Appendix.

141 **Notation.** We denote by \mathbb{R} and \mathbb{R}_+ , respectively, the set of real numbers and real
 142 nonnegative numbers, and by \mathbb{R}^m the m -dimensional Euclidean real coordinate space.
 143 We use normal (lower-case or upper-case) letters for scalars, lower-case boldface letters
 144 for vectors, and upper case boldface letters for matrices. Further, we denote by: a_i or
 145 $[\mathbf{a}]_i$, as appropriate, the i -th element of vector \mathbf{a} ; \mathbf{A}_{ij} or $[\mathbf{A}]_{ij}$, as appropriate, the entry
 146 in the i -th row and j -th column of a matrix \mathbf{A} ; \mathbf{A}^\top the transpose of a matrix \mathbf{A} ; and
 147 $\text{trace}(\mathbf{A})$ the sum of diagonal elements of \mathbf{A} . Further, we use either $\mathbf{a}^\top \mathbf{b}$ or $\langle \mathbf{a}, \mathbf{b} \rangle$ for
 148 the inner product of vectors \mathbf{a} and \mathbf{b} . Next, we let \mathbf{I} and $\mathbf{0}$ be, respectively, the identity
 149 matrix and the zero matrix; $\|\cdot\| = \|\cdot\|_2$ the Euclidean (respectively, spectral) norm
 150 of its vector (respectively, matrix) argument; $\phi'(w)$ the first derivative evaluated at w
 151 of a function $\phi : \mathbb{R} \rightarrow \mathbb{R}$; $\nabla h(\mathbf{w})$ and $\nabla^2 h(\mathbf{w})$ the gradient and Hessian, respectively,
 152 evaluated at \mathbf{w} of a function $h : \mathbb{R}^m \rightarrow \mathbb{R}$; $\mathbb{P}(\mathcal{A})$ and $\mathbb{E}[u]$ the probability of an event \mathcal{A}
 153 and expectation of a random variable u , respectively; and by $\text{sign}(a)$ the sign function,
 154 i.e., $\text{sign}(a) = 1$, for $a > 0$, $\text{sign}(a) = -1$, for $a < 0$, and $\text{sign}(0) = 0$. Finally, for
 155 two positive sequences η_n and χ_n , we have: $\eta_n = O(\chi_n)$ if $\limsup_{n \rightarrow \infty} \frac{\eta_n}{\chi_n} < \infty$;
 156 $\eta_n = \Omega(\chi_n)$ if $\liminf_{n \rightarrow \infty} \frac{\eta_n}{\chi_n} > 0$; and $\eta_n = \Theta(\chi_n)$ if $\eta_n = O(\chi_n)$ and $\eta_n = \Omega(\chi_n)$.

157 **2. Problem Model and the nonlinear SGD Framework.** We consider the
 158 following unconstrained problem:

$$159 \quad (2.1) \quad \text{minimize } f(\mathbf{x}),$$

161 where $f : \mathbb{R}^d \mapsto \mathbb{R}$ is a convex function.

162 We make the following standard assumption.

163 **ASSUMPTION 1.** *Function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is strongly convex with strong convexity*
 164 *parameter $\mu > 0$, and it has Lipschitz continuous gradient with Lipschitz constant $L \geq$*
 165 *μ .*

166 For asymptotic results (see ahead [Theorems 3.1](#) and [3.3](#)), we will also require the
 167 following assumption.

168 **ASSUMPTION 2.** *Function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is twice continuously differentiable.*

169 Under [Assumption 1](#), problem (2.1) has a unique solution, which we denote by
 170 $\mathbf{x}^* \in \mathbb{R}^d$.

171 In machine learning settings, f can correspond to the risk function, i.e.,

$$172 \quad (2.2) \quad f(\mathbf{x}) = \mathbb{E}_{\mathbf{d} \sim P} [\ell(\mathbf{x}; \mathbf{d})] + \mathcal{R}(\mathbf{x}).$$

173 Here, P is the (unknown) distribution from which the data samples $\mathbf{d} \in \mathbb{R}^q$ are
 174 drawn; $\ell(\cdot; \cdot)$ is a loss function that is smooth and convex in its first argument for any
 175 fixed value of the second argument; and $\mathcal{R} : \mathbb{R}^d \mapsto \mathbb{R}$ is a smooth strongly convex
 176 regularizer. Similarly, f can be empirical risk, i.e., $f(\mathbf{x}) = \frac{1}{n} \left(\sum_{j=1}^n \ell(\mathbf{x}; \mathbf{d}_j) \right) + \mathcal{R}(\mathbf{x})$,
 177 where \mathbf{d}_j , $j = 1, \dots, n$, is the set of training data points. Several machine learning
 178 models fall within the described framework under [Assumptions 1–2](#), including, e.g.,
 179 ℓ_2 -regularized quadratic and logistic losses.

180 We introduce a general framework for *nonlinear* SGD methods to solve prob-
 181 lem (1); an algorithm within the framework takes the following form:

$$182 \quad (2.3) \quad \mathbf{x}^{t+1} = \mathbf{x}^t - \alpha_t \Psi(\nabla f(\mathbf{x}^t) + \boldsymbol{\nu}^t).$$

183 Here, \mathbf{x}^t denotes the solution estimate at iteration t , $t = 0, 1, \dots$; $\Psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ is
 184 a general nonlinear map; $\alpha_t > 0$ is the employed step size; $\boldsymbol{\nu}^t \in \mathbb{R}^d$ is a zero-mean
 185 gradient noise; and \mathbf{x}^0 is an arbitrary deterministic point in \mathbb{R}^d .

186 We will specify further ahead the assumptions that we make on the step size α_t ,
 187 the map Ψ and the noise $\boldsymbol{\nu}^t$. Some examples of commonly used maps Ψ that fall
 188 within our framework are the following:

- 189 1. Sign gradient: $[\Psi(\mathbf{w})]_i = \text{sign}(w_i)$, $i = 1, \dots, d$;
- 190 2. Component-wise clipping: $[\Psi(\mathbf{w})]_i = w_i$, for $|w_i| \leq m$; $[\Psi(\mathbf{w})]_i = m$, for
 191 $w_i > m$, and $[\Psi(\mathbf{w})]_i = -m$, for $w_i < -m$, for some constant $m > 0$.
- 192 3. Component-wise quantization: for each $i = 1, \dots, d$, we let $[\Psi(\mathbf{w})]_i = r_j$, for
 193 $w_i \in (q_{j-1}, q_j]$, $j = 1, \dots, J$, where $-\infty = q_0 < q_1 < \dots < q_J = +\infty$, J is a
 194 positive integer, and the r_j 's and q_j 's are chosen such that each component
 195 nonlinearity is an odd function, i.e., $[\Psi(\mathbf{w})]_i = -[\Psi(-\mathbf{w})]_i$, for each i and for
 196 each \mathbf{w} ;
- 197 4. Normalized gradient: $\Psi(\mathbf{w}) = \frac{\mathbf{w}}{\|\mathbf{w}\|}$, for $\mathbf{w} \neq 0$, and $\Psi(0) = 0$;
- 198 5. Clipped gradient: $\Psi(\mathbf{w}) = \mathbf{w}$, for $\|\mathbf{w}\| \leq M$, and $\Psi(\mathbf{w}) = \frac{\mathbf{w}}{\|\mathbf{w}\|} M$, for
 199 $\|\mathbf{w}\| > M$, for some constant $M > 0$.

200 Other nonlinearity choices are also introduced ahead (see [Section 7](#)).

201 We next discuss the various possible sources of the gradient noise $\boldsymbol{\nu}^t$. First, the
 202 noise may arise due to utilizing a search direction with respect to a data sample. That
 203 is, a common search direction in machine learning algorithms is the gradient of the loss
 204 with respect to a single data point \mathbf{d}_i : $\mathbf{g}_i(\mathbf{x}) = \nabla \ell(\mathbf{x}; \mathbf{d}_i) + \nabla \mathcal{R}(\mathbf{x})$. In case of the risk
 205 function (2.2), \mathbf{d}_i is drawn from distribution P ; in case of the empirical risk, \mathbf{d}_i
 206 can be, e.g., drawn uniformly at random from the set of data points \mathbf{d}_j , $j = 1, \dots, n$, with
 207 repetition along iterations. In both cases, the corresponding gradient noise equals
 208 $\boldsymbol{\nu} = \mathbf{g}_i(\mathbf{x}) - \nabla f(\mathbf{x})$. Several recent studies indicate that noise $\boldsymbol{\nu}$ exhibits heavy tails
 209 on many real data sets, e.g. [32, 17, 37]. (See also [Section 7](#)).

210 We also comment on other possible sources of gradient noise. The noise may
 211 be added on purpose to the gradient $\nabla f(\mathbf{x})$ for improving privacy of an SGD-based
 212 learning process, e.g., [29]. Also, the noise $\boldsymbol{\nu}^t$ may model random computational
 213 perturbations or inexact calculations in evaluating a gradient $\nabla f(\mathbf{x})$.

214 **3. Main results: Component-wise Nonlinearities.** Section 3 provides anal-
 215 ysis of the nonlinear SGD method for component-wise nonlinearities. That is, we
 216 consider here maps $\Psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ of the form $\Psi(w_1, \dots, w_d) = (\Psi(w_1), \dots, \Psi(w_d))^T$,
 217 for any $\mathbf{w} \in \mathbb{R}^d$, where (somewhat abusing notation) we denote by $\Psi : \mathbb{R} \mapsto \mathbb{R}$ the
 218 component-wise nonlinearity. In this setting, we establish for (2.3) almost sure con-
 219 vergence and evaluate the MSE convergence rate and the asymptotic covariance of
 220 the method. In more detail, we consider a probability space (Ω, \mathcal{F}, P) , where $\omega \in \Omega$ is
 221 a canonical element. For each $t = 0, 1, \dots$, $\boldsymbol{\nu}^t : \Omega \mapsto \mathbb{R}^d$ is a random vector defined on
 222 (Ω, \mathcal{F}, P) . We also denote by \mathcal{F}_t , $t = 0, 1, \dots$, the σ -algebra generated by random vec-
 223 tors $\{\boldsymbol{\nu}^s\}$, $s = 0, \dots, t$. Clearly, in view of (2.3), \mathbf{x}^{t+1} is measurable with respect to \mathcal{F}_t ,
 224 $t = 0, 1, \dots$. We make the following assumptions; they follow the noise and nonlinearity
 225 framework similar to [30].

226 **ASSUMPTION 3 (Gradient noise).** *For the gradient noise random vector sequence*
 227 *$\{\boldsymbol{\nu}^t\}$ in (2.3), $t = 0, 1, \dots$, $\boldsymbol{\nu}^t \in \mathbb{R}^d$, we assume the following.*

- 228 1. *The sequence of random vectors $\{\boldsymbol{\nu}^t\}$ is independent identically distributed*

²Similar considerations hold for a loss with respect to a mini-batch of data points; this discussion is abstracted for simplicity.

- 229 (i.i.d.) Also, random variables ν_i^t are mutually independent across $i = 1, \dots, d$;
 230 2. Each component ν_i^t , $i = 1, \dots, d$, of vector $\boldsymbol{\nu}^t = (\nu_1^t, \dots, \nu_d^t)^\top$ has a probability
 231 density function $p(u)$, $p : \mathbb{R} \mapsto \mathbb{R}_+$.
 232 3. The pdf p is symmetric, i.e., $p(u) = p(-u)$, for any $u \in \mathbb{R}$ with $\int |u|p(u)du <$
 233 $+\infty$, and $p(u) > 0$ for $|u| \leq c_p$, for some constant $c_p > 0$.

234 Note that Assumption 3 implies that $\boldsymbol{\nu}^t$ is zero-mean, for all t , and that $\boldsymbol{\nu}^t$ and \mathbf{x}^t
 235 are mutually independent, for all t . For a class of unbounded nonlinearities Ψ that
 236 obey Assumption 6 ahead, we will additionally require the following.

237 ASSUMPTION 4. The gradient noise variance $\sigma_\nu^2 = \int_{-\infty}^{+\infty} u^2 p(u) du < +\infty$.

238 Assumption 3 requires that the noise vector is i.i.d. across its components $i = 1, \dots, d$
 239 which may be restrictive in certain scenarios. For the global MSE analysis, these
 240 assumptions can be relaxed; see ahead the remark after Theorem 3.2 and Appendix C.

241 Regarding noise pdf $p(u)$, except for strictly positive values in the vicinity of
 242 zero (a very mild assumption), we require that the noise pdf is symmetric. Examples
 243 of the distributions that satisfy Assumption 3 include, e.g., a Gaussian zero-mean pdf
 244 or a Laplace zero-mean pdf with strictly positive variances, and heavy-tail zero-mean
 245 symmetric α -stable distributions [3].³ On the other hand, $p(u)$ may not be symmetric
 246 if, e.g., it is a mixture of some standard distributions. For example, consider random
 247 variable ν that is sampled from $\mathbb{N}(-m_1, \sigma^2)$ with probability $p = \frac{m_2}{m_1 + m_2}$ and it is
 248 sampled from $\mathbb{N}(m_2, \sigma^2)$ with probability $1 - p$, for some $m_1 \neq m_2$, $m_1, m_2 > 0$, and
 249 $\sigma > 0$. Then, clearly, ν is zero-mean but does not have a symmetric pdf.

250 ASSUMPTION 5 (Nonlinearity Ψ). Function $\Psi : \mathbb{R} \mapsto \mathbb{R}$ is a continuous (except
 251 possibly on a point set with Lebesgue measure of zero), monotonically non-decreasing
 252 and odd function, i.e., $\Psi(-w) = -\Psi(w)$, for any $w \in \mathbb{R}$. Moreover, Ψ is piece-wise
 253 differentiable. Finally, Ψ is either discontinuous at zero, or $\Psi(u)$ is strictly increasing
 254 for $u \in (-c_\Psi, c_\Psi)$, for some $c_\Psi > 0$.

255 In addition, we impose one of the Assumptions 6 or 7 below.

256 ASSUMPTION 6. $|\Psi(w)| \leq C_1 (1 + |w|)$, for any $w \in \mathbb{R}$, for some constant $C_1 > 0$.

257 ASSUMPTION 7. $|\Psi(w)| \leq C_2$, for some constant $C_2 > 0$.

258 Assumption 3 and Assumption 5 are imposed throughout the paper. Assumption 4
 259 is imposed when Assumption 6 holds, i.e., for the nonlinearities Ψ that can have un-
 260 bounded outputs. When Assumption 7 is imposed, then Assumption 4 is not required.

261 Note that, provided that Assumption 7 holds, we require only a finite first moment
 262 of the gradient noise, while the moments of α -order, $\alpha > 1$, may be infinite, hence
 263 allowing for heavy-tail noise distributions. For example, the gradient noise variance
 264 can be infinite. Assumption 5 holds for several interesting component-wise nonlinear-
 265 ities, like, e.g., the sign gradient, component-wise clipping, and quantization schemes
 266 introduced in Section 2. Note also that Assumption 5 encompasses a broad range of
 267 component-wise nonlinearities, beyond the examples in Section 2. (For example, see
 268 Section 7 for the tanh and a bi-level quantization nonlinearity.)

269 Let us define function $\phi : \mathbb{R} \mapsto \mathbb{R}$, as follows. For a fixed (deterministic) point
 270 $w \in \mathbb{R}$, $\phi(w)$ is defined by:

$$271 \quad (3.1) \quad \phi(w) = \mathbb{E} [\Psi(w + \nu_1^0)] = \int \Psi(w + u)p(u)du,$$

³A random variable Z has a symmetric α -stable zero-mean distribution with scale parameter $\sigma > 0$ if its characteristic function takes the form: $\mathbb{E}[\exp(iuZ)] = \exp(-\sigma^\alpha |u|^\alpha)$, $u \in \mathbb{R}$, $\alpha \in [0, 2]$.

272 where the expectation is taken with respect to the distribution of a single entry of
 273 the gradient noise at any iteration, i.e., with respect to pdf $p(u)$. Intuitively, the
 274 nonlinearity ϕ is a convolution-like transformation of the nonlinearity Ψ , where the
 275 convolution is taken with respect to the gradient noise pdf $p(u)$. As we will see
 276 ahead, the nonlinearity ϕ plays an effective role in determining the performance
 277 of algorithm (2.3). We now state the main results on (2.3) with component-wise
 278 nonlinearities, including the results on a.s. convergence, MSE rate, and asymptotic
 279 normality. We start with the following Theorem that establishes a.s. convergence.

280 **THEOREM 3.1** (Almost sure convergence: Component-wise nonlinearity). *Con-*
 281 *sider algorithm (2.3) for solving optimization problem (2.1), and let Assumptions*
 282 *1, 2, 3, 5, and 7 hold. Further, let the positive step-size sequence $\{\alpha_t\}$ be square*
 283 *summable, non-summable: $\sum \alpha_t = +\infty$; $\sum \alpha_t^2 < +\infty$. Then, the sequence of iterates*
 284 *$\{\mathbf{x}^t\}$ generated by algorithm (2.3) converges almost surely to the solution \mathbf{x}^* of the*
 285 *optimization problem (2.1). Moreover, the result holds if Assumption 7 is replaced*
 286 *with Assumption 6, and Assumption 4 is additionally imposed.*

287 **Theorem 3.1** establishes a.s. convergence of the nonlinear SGD scheme (2.3)
 288 under a general setting for the component-wise nonlinearities and gradient noise. For
 289 example, provided that the output of the nonlinearity Ψ is bounded, algorithm (2.3)
 290 converges even when the gradient noise may not have a finite α -moment, for any $\alpha > 1$.
 291 (Hence it may have an infinite variance). In contrast, as shown in Appendix B, the
 292 linear SGD (algorithm (2.3) with Ψ being the identity function) generates a sequence
 293 of solution estimates with infinite variances, provided that the variance of $p(u)$ is
 294 infinite.

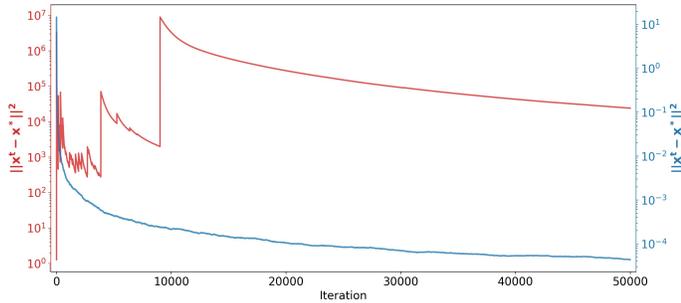


Fig. 1: Illustration of **Theorem 3.1**: estimated MSE versus iteration counter for the nonlinear SGD in (2.3) with component-wise sign nonlinearity (blue line) and the linear SGD (red line).

295 **Example 3.1.** **Figure 1** illustrates **Theorem 3.1** with a simulation example. We
 296 consider a strongly convex quadratic function $f : \mathbb{R}^d \mapsto \mathbb{R}$, $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b}^\top \mathbf{x}$,
 297 where $\mathbf{A} \in \mathbb{R}^{d \times d}$ is a (symmetric) positive definite matrix, $d = 16$, and quantities
 298 \mathbf{A}, \mathbf{b} are generated at random. We consider algorithm (2.3) with the component-wise
 299 sign nonlinearity and the linear SGD. The gradient noise has a heavy-tailed pdf given
 300 by:

301 (3.2)
$$p(u) = \frac{\alpha - 1}{2(1 + |u|)^\alpha},$$

for $u \in \mathbb{R}$ and $\alpha > 2$. Note that the distribution (3.2) does not have a finite $\alpha - 1$ moment and has finite moments of r -th order for $r < \alpha - 1$. We set in simulation $\alpha = 2.05$. Note that, in this case, the gradient noise has infinite variance. We initialize both the linear and nonlinear algorithm with $\mathbf{x}^0 = 0$, and we let step size $\alpha_t = \frac{1}{t+1}$. Figure 1 shows an estimate of MSE, i.e., of the quantity $\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|^2]$, obtained by averaging results from 100 sample paths. The red line corresponds to the linear SGD, while the blue line corresponds to the nonlinear SGD with the component-wise sign nonlinearity. As predicted by Theorem 3.1, the nonlinear SGD drives the MSE to zero, while the linear SGD does not seem to provide a meaningful solution estimate sequence.

We next establish the mean square error (MSE) convergence rate of algorithm (2.3).

THEOREM 3.2 (MSE convergence: Component-wise nonlinearity). *Consider algorithm (2.3) for solving optimization problem (2.1), and let Assumptions 1, 3, 5, and 7 hold. Further, let the step-size sequence $\{\alpha_t\}$ be $\alpha_t = a/(t+1)^\delta$, $a > 0$, $\delta \in (0.5, 1)$. Then, for the sequence of iterates $\{\mathbf{x}^t\}$ generated by algorithm (2.3), it holds that $\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|^2] = O(1/t^\zeta)$, or equivalently, $\mathbb{E}[f(\mathbf{x}^t) - f^*] = O(1/t^\zeta)$. Here, $\zeta < 1$ is any positive number such that $\zeta < \min\left(2\delta - 1, \frac{a(1-\delta)\xi\phi'(0)\mu}{L(aC_2\sqrt{a} + \|\mathbf{x}^0 - \mathbf{x}^*\|)}\right)$, and constant $\xi > 0$ is such that $\phi(a) \geq \frac{\phi'(0)}{2}a$, for any $a \in [0, \xi]$. Furthermore, let Assumptions 1, 3, 5, and 6, and 4 hold, let $\alpha_t = \frac{a}{(t+1)^\delta}$, $\delta \in (0.5, 1]$, and assume that $\inf_{a \neq 0} \frac{|\Psi(a)|}{|a|} > 0$. Then, there holds that $\mathbb{E}[\|\mathbf{x}^t - \mathbf{x}^*\|^2] = O(1/t^\delta)$, or equivalently, $\mathbb{E}[f(\mathbf{x}^t) - f^*] = O(1/t^\delta)$. In particular, for $\delta = 1$, we obtain the $O(1/t)$ MSE rate.*

Remark. The MSE convergence $O(1/t^{\zeta'})$, for some $\zeta' \in (0, 1)$, continues to hold under the same set of assumptions as in Theorem 3.2 but with a relaxed version of Assumption 3, where we no longer require that the gradient noise vector has mutually independent components. More precisely, we allow for an i.i.d. noise vector sequence $\{\boldsymbol{\nu}^t\}$, $\boldsymbol{\nu}^t \in \mathbb{R}^d$, that has a symmetric joint pdf $p : \mathbb{R}^d \mapsto \mathbb{R}$, $p(\mathbf{u}) = p(-\mathbf{u})$, for any $\mathbf{u} \in \mathbb{R}^d$, that is strictly positive for $\|\mathbf{u}\| \leq u_0$, for some $u_0 > 0$. In that case, effectively, the role of function ϕ in Theorem 3.2 is replaced by functions $w \mapsto \phi_i(w)$, $w \in \mathbb{R}$, $i = 1, \dots, d$, where $\phi_i(w) = \int \Psi(w+u)p_i(u)du$, and $p_i : \mathbb{R} \mapsto \mathbb{R}$ is the marginal pdf of the i -th component associated with the joint pdf $p : \mathbb{R}^d \mapsto \mathbb{R}$. (See Appendix C.)

For the bounded nonlinearity case (e.g., sign gradient, component-wise clipping, quantization nonlinearity) and the heavy-tail noise (only the first noise moment assumed to be finite), the nonlinear SGD (2.3) achieves a global sublinear MSE rate $O(1/t^\zeta)$, $\zeta \in (0, 1)$. On the other hand, for the finite variance case and an unbounded nonlinearity, the nonlinear SGD (2.3) achieves a global MSE rate $O(1/t)$ provided that $\inf_{w \neq 0} \frac{|\Psi(w)|}{|w|} > 0$. This is the best achievable rate and equal to that of the linear SGD in the same setting. Furthermore, by Theorem 3.3 ahead, the nonlinear SGD (2.3) with bounded outputs under the heavy-tail noise achieves *locally*, in the weak convergence sense, the faster $O(1/t)$ rate. This is again in the setting where the linear SGD fails.

Example 3.2. We next illustrate the value ζ in Theorem 3.2 on the family of heavy-tailed pdfs given in (3.2). To be specific, consider the sign nonlinearity $\Psi(w) = \text{sign}(w)$. Then, it is easy to show that: $\phi(w) = 2 \int_0^w p(u)du$, $\phi'(0) = 2p(0)$, $\xi \geq 2^{1/\alpha} - 1 \approx \frac{1}{\alpha}$. Using the above calculations, we can see that, for a large a , ζ can be approximated as $\min\{2\delta - 1, \frac{\mu}{L} \frac{1-\delta}{\sqrt{a}} \frac{\alpha-1}{\alpha}\}$.

We also compare the rate ζ with the analysis in [37] that is closest to our setting

349 with respect to existing work. Modulo the differences in the assumptions of the
 350 assumed settings here and in [37], the rate in [37], when adapted to the noise pdf
 351 in Example 3.1, reads as follows: $\frac{2(r-1)}{r}$, where r is any number such that $r \leq$
 352 $\min\{\alpha - 1, 2\}$. When compared with ζ , the rate in [37] is clearly better for α above a
 353 threshold. However, as α decreases and approaches the value 2, the rate achieved here
 354 stays bounded away from zero and approaches the quantity: $\min\left\{2\delta - 1, \frac{1}{2} \frac{\mu}{L} \frac{1-\delta}{\sqrt{d}}\right\}$.
 355 In contrast, the rate in in [37] approaches zero as α approaches 2. ⁴

356 **Example 3.3.** We continue to assume the noise pdf in (3.2), but here we consider
 357 the component-wise clipping nonlinearity Ψ with saturation value m . For simplicity,
 358 we take $m > 1$, while similar bounds can be obtained for $m \leq 1$ as well. It can be
 359 shown that the rate ζ can be estimated as (see Appendix E):

$$360 \quad (3.3) \quad \min\left\{2\delta - 1, \frac{\mu}{L\sqrt{d}} \frac{(1-\delta)(m-1)(1-(m+1)^{-\alpha})}{m}\right\}.$$

361 The above α -dependent estimate can be replaced with a more conservative rate that
 362 holds for any $\alpha > 2$: $\min\left\{2\delta - 1, \frac{\mu}{L\sqrt{d}} \frac{(1-\delta)(m-1)(1-(m+1)^{-2})}{m}\right\}$. We again compare
 363 the rate achieved by the proposed method with the rate from [37] that equals: $\frac{2(r-1)}{r}$,
 364 $r < \min\{\alpha - 1, 2\}$. We can see that the rate in [37] is better than (3.3) for α above
 365 a threshold. On the other hand, when α decreases to 2, the rate of [37] approaches
 366 zero, while (3.3) becomes better and stays bounded away from zero.

367 We next establish asymptotic normality of (2.3).

368 **THEOREM 3.3** (Asymptotic normality: Component-wise nonlinearity). *Consider*
 369 *algorithm (2.3) for solving optimization problem (2.1), and let Assumptions 1, 2,*
 370 *3, 5, and 7 hold. Further, let the step-size sequence $\{\alpha_t\}$ equal: $\alpha_t = a/(t+1)$,*
 371 *$t = 0, 1, \dots$, with parameter $a > \frac{1}{2\phi'(0)\mu}$. Then, the sequence of iterates $\{\mathbf{x}^t\}$ generated*
 372 *by algorithm (2.3) is asymptotically normal, and there holds:*

$$373 \quad (3.4) \quad \sqrt{t+1}(\mathbf{x}^t - \mathbf{x}^*) \xrightarrow{d} \mathbb{N}(0, \mathcal{S}),$$

374 where \xrightarrow{d} designates convergence in distribution. The asymptotic covariance \mathcal{S} of the
 375 multivariate normal distribution $\mathbb{N}(0, \mathcal{S})$ is given by:

$$376 \quad \mathcal{S} = a^2 \int_{\nu=0}^{\infty} e^{\nu\Sigma} \mathcal{S}_0 e^{\nu\Sigma} d\nu = a^2 \sigma_{\Psi}^2 [2a\phi'(0)\nabla^2 f(x^*) - \mathbf{I}]^{-1},$$

377 where:

$$379 \quad (3.5) \quad \mathcal{S}_0 = \sigma_{\Psi}^2 \mathbf{I}, \quad \sigma_{\Psi}^2 = \int |\Psi(v)|^2 p(v) dv, \quad \Sigma = \frac{1}{2} \mathbf{I} - a \phi'(a) \nabla^2 f(\mathbf{x}^*).$$

380 Moreover, the same result holds when Assumption 7 is replaced with Assumption 6,
 381 and Assumption 4 is additionally imposed.

⁴It is worth noting that reference [37] establishes certain tightness results on the rate achieved therein, by providing a “hard” problem example where the mean squared error after t iterations is $\Omega(1/t^{\frac{2(r-1)}{r}})$. However, this does not contradict our results due to the different sets of Assumptions made here and in [37]. Most notably, [37] assumes bounded moments of gradients and allow for dependence between the current point \mathbf{x}^t and the gradient noise ν^t . In fact, the “hard example” construction in the proof of Theorem 5 in [37] constructs ν^t as an explicit function of \mathbf{x}^t .

382 **Theorem 3.3** establishes asymptotic normality of (2.3) and, moreover, it gives an
 383 exact expression for the asymptotic covariance \mathcal{S} in (3.3), that basically corresponds
 384 to the constant in the $1/t$ variance decay near the solution. The asymptotic covariance
 385 value (3.3) reveals an interesting tradeoff with respect to the effect of the nonlinearity
 386 Ψ . We provide some insights into the tradeoff through examples below.

387 **Example 3.4** Figure 2 illustrates **Theorem 3.3** for the nonlinear SGD in (2.3)
 388 with component-wise sign nonlinearity and the same simulation setting used for the
 389 numerical illustration of **Theorem 3.1** and step-size $\alpha_t = \frac{10}{t+1}$. The red line plots quantity
 390 $\frac{t}{d} \|\mathbf{x}^t - \mathbf{x}^*\|^2$ estimated through 100 sample path runs. This quantity estimates
 391 the constant in the $1/t$ per-entry asymptotic variance decay, i.e., it is a numerical estimate
 392 of the per-entry asymptotic variance $\frac{\text{trace}(\mathcal{S})}{d}$, where \mathcal{S} is given in **Theorem 3.3**.
 393 The blue horizontal line marks the value $\frac{\text{trace}(\mathcal{S})}{d}$. We can see that the simulation
 394 matches well the theory.

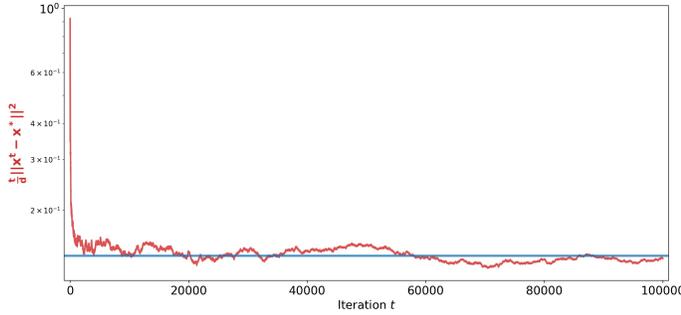


Fig. 2: Illustration of **Theorem 3.3**: Monte Carlo estimate of per-entry asymptotic variance (red line) and the theoretical per-entry asymptotic variance in **Theorem 3.3** (blue line).

395 **Example 3.5.** We compare the linear SGD and the nonlinear SGD with
 396 component-wise clipping. For illustration and simplification of calculations, we consider
 397 the special case when $\nabla^2 f(\mathbf{x}^*)$ is a symmetric matrix with all eigenvalues equal
 398 to one. Then, it is straightforward to show that the per-entry asymptotic variance for
 399 the best choice of parameter a over the admissible set of values equals:

$$400 \quad (3.6) \quad \inf_{a > \frac{1}{2\phi'(0)}} \text{trace}(\mathcal{S}) = \frac{\sigma_{\Psi}^2}{(\phi'(0))^2}.$$

401 Here, for the linear SGD i.e., when $\Psi(a) = a$, we have that $\sigma_{\Psi}^2 = \int a^2 p(a) da$ equals
 402 the gradient noise (per component) variance σ_v^2 , and $\phi'(0) = 1$, and so (3.6) equals σ_v^2 .
 403 Now, consider the coordinate-wise clipping, with $\Psi(a) = a$ for $|a| \leq m$ and $\Psi(a) =$
 404 $\text{sign}(a)m$, for $|a| > m$, for some $m > 0$. Then, we have: $\sigma_{\Psi}^2 = m^2 - 2 \int_0^m (m^2 -$
 405 $v^2)p(v)dv$, and $\phi'(0) = 2 \int_0^m p(v)dv$. (See Appendix F for the derivation.) Note that
 406 the case $m \rightarrow \infty$ corresponds to the linear SGD case. Consider now the tradeoff with
 407 respect to the choice of m . Clearly, taking a smaller m has a positive effect on the
 408 numerator in (3.6) (it suppresses the noise effect). On the other hand, reducing m
 409 has a negative effect on the denominator in (3.6); that is, it reduces the value $\phi'(0)$
 410 – intuitively, it “lowers the quality” of the search direction utilized with (2.3). One
 411 needs to choose the nonlinearity, i.e., the parameter m , optimally, to strike the best

412 balance here. Clearly, for larger gradient noise σ_ν^2 , we should pick a smaller value
 413 of m . Note also that, when σ_ν^2 is infinite, the linear SGD has an infinite asymptotic
 414 variance in (3.6), while the nonlinear SGD with any $m \in (0, \infty)$ has a finite asymptotic
 415 variance.

416 **Example 3.6.** We continue to assume the simplified setting when the per-entry
 417 asymptotic variance equals (3.6). We consider the sign gradient nonlinearity and
 418 the class of heavy-tail gradient noise distributions in (3.2). It can be shown that
 419 here: $\sigma_\Psi^2 = 1$; $\sigma_\nu^2 = \frac{2}{(\alpha-3)(\alpha-2)}$, for $\alpha > 3$ and $\sigma_\nu^2 = \infty$, else; and $\phi'(0) = \alpha - 1$.
 420 (See Appendix G.) Therefore, for the sign gradient, the best achievable per entry
 421 asymptotic variance equals $\frac{1}{(\alpha-1)^2}$, while for the linear SGD it equals $\frac{2}{(\alpha-2)(\alpha-3)}$ for
 422 $\alpha > 3$, and is infinite for $\alpha \in (2, 3]$. Hence, we can see for the considered example that
 423 the sign gradient outperforms the linear SGD for any $\alpha > 2$, and the gap becomes
 424 larger as α gets smaller.

425 **Example 3.7.** We still consider the simplified setting of (3.6). If the noise pdf
 426 $p(u)$ is known, then, following [30], we can find a globally optimal nonlinearity that
 427 minimizes (3.6) that takes the form: $\Psi(a) = -\frac{d}{da} \ln(p(a))$. The corresponding optimal
 428 asymptotic variance equals the Fisher information associated with the pdf $p(u)$.

429 **4. Main results: Joint Nonlinearities.** We now consider algorithm (2.3) for
 430 a nonlinearity $\Psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ that cannot be decoupled into (equal) component wise
 431 nonlinearities $\Psi : \mathbb{R} \mapsto \mathbb{R}$, as it was possible before. More precisely, we make the
 432 following assumptions on the gradient noise ν^t and the nonlinear map $\Psi : \mathbb{R}^d \mapsto \mathbb{R}^d$.
 433 Recall also filtration \mathcal{F}_t in Section 3.

434 **ASSUMPTION 8.** [Gradient noise] For the gradient noise sequence $\{\nu^t\}$, we as-
 435 *sume the following:*

- 436 1. The sequence of random vectors $\{\nu^t\}$ is i.i.d. Moreover, ν^t has a joint sym-
 437 metric pdf $p(\mathbf{u})$, $p : \mathbb{R}^d \mapsto \mathbb{R}$, i.e., $p(\mathbf{u}) = p(-\mathbf{u})$, for any $\mathbf{u} \in \mathbb{R}^d$ with
 438 $\int \|\mathbf{u}\| p(\mathbf{u}) d\mathbf{u} < \infty$;
- 439 2. There exists a positive constant B_0 such that, for any $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{x} \neq 0$, for any
 440 $A \in (0, 1]$, there exists $\lambda = \lambda(A) > 0$, such that $\int_{\mathcal{J}_A} p(\mathbf{u}) d\mathbf{u} > \lambda(A)$, where
 441 $\mathcal{J}_A = \{\mathbf{u} \in \mathbb{R}^d : \frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\| \|\mathbf{x}\|} \in [0, A], \|\mathbf{u}\| \leq B_0\}$.⁵

442 **Assumption 8** allows for a heavy-tailed noise vector whose components can be
 443 mutually dependent. Condition 2. in **Assumption 8** is mild; it says that the joint
 444 pdf $p(\mathbf{u})$ is “non-degenerate” in the sense that, along each “direction” (determined
 445 by arbitrary nonzero vector \mathbf{x}), the intersection of the set $\{\frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\| \|\mathbf{x}\|} \in [0, A]\}$ and the
 446 ball $\{\|\mathbf{u}\| \leq B_0\}$ consumes a positive mass of the joint pdf $p(\mathbf{u})$.

447 We make the following assumption on the joint nonlinearity.

448 **ASSUMPTION 9 (Nonlinearity Ψ).** The nonlinear map $\Psi : \mathbb{R}^d \mapsto \mathbb{R}^d$ takes the
 449 following form: $\Psi(\mathbf{w}) = \mathbf{w} \mathcal{N}(\|\mathbf{w}\|)$, where function $\mathcal{N} : \mathbb{R}_+ \mapsto \mathbb{R}_+$ satisfies the
 450 following: \mathcal{N} is non-increasing and continuous except possibly on a point set with
 451 Lebesgue measure of zero with $\mathcal{N}(q) > 0$, for any $q > 0$. The function $q\mathcal{N}(q)$ is
 452 non-decreasing.

453 In addition, we assume that either **Assumption 10** or **Assumption 11** holds.

454 **ASSUMPTION 10.** $\|\Psi(\mathbf{w})\| \leq C'_2$, for any $\mathbf{w} \in \mathbb{R}^d$, for some $C'_2 > 0$.

⁵The integration set \mathcal{J}_A also includes the point $\mathbf{u} = 0$. In other words, for compact notation here and throughout the paper, we write $\frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\| \|\mathbf{x}\|} \in [0, A]$ instead of $0 \leq \mathbf{u}^\top \mathbf{x} \leq A \|\mathbf{u}\| \|\mathbf{x}\|$.

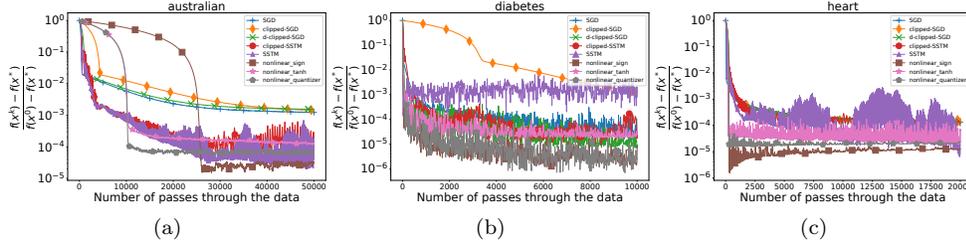


Fig. 3: Comparison of the optimization algorithms across different datasets

455 ASSUMPTION 11. $\|\Psi(\mathbf{w})\| \leq C'_1(1 + \|\mathbf{w}\|)$, for any $\mathbf{w} \in \mathbb{R}^d$, for some $C'_1 > 0$

456 There are many nonlinearities that satisfy the above Assumptions, including the nor-
 457 malized gradient and the clipped gradient discussed in Section 2. If Assumption 11
 458 holds, then we additionally require the following.

459 ASSUMPTION 12. *There holds: $\int \|\mathbf{u}\|^2 p(\mathbf{u}) d\mathbf{u} < \infty$.*

460 For asymptotic normality in the joint nonlinearity case, we additionally impose
 461 the following.

462 ASSUMPTION 13. *Function $\mathcal{N} : \mathbb{R}_+ \mapsto \mathbb{R}$ is differentiable for any positive argu-
 463 ment, i.e., $\mathcal{N}'(a)$ exists for any $a > 0$. Furthermore, $\sup_{a>0} \mathcal{N}(a) < +\infty$.*

464 We first state Theorem 4.1 and Theorem 4.2 on the a.s. convergence and the MSE
 465 rate of algorithm (2.3), respectively; we then illustrate the results with examples.

466 THEOREM 4.1 (A.s. convergence: Joint nonlinearity). *Consider algorithm (2.3)
 467 for solving optimization problem (2.1), and let Assumptions 1, 2, 8, 9, and 10 hold.
 468 Further, let the step-size sequence $\{\alpha_t\}$ be square-summable, non-summable. Then,
 469 for the sequence of iterates $\{\mathbf{x}^t\}$ generated by algorithm (2.3), it holds that $\mathbf{x}^t \rightarrow$
 470 \mathbf{x}^* , a.s. Moreover, the result continues to hold if Assumption 10 is replaced with
 471 Assumption 11, and Assumption 12 is additionally imposed.*

472 We now state our MSE rate result for the joint nonlinearity case.

473 THEOREM 4.2 (MSE convergence rate: Joint nonlinearity). *Consider algorithm
 474 (2.3) for solving optimization problem (2.1), and let Assumptions 1, 8, 9, and 10 hold.
 475 Further, let the step-size sequence $\{\alpha_t\}$ be $\alpha_t = a/(t+1)$, $a > 0$, $\delta \in (0.5, 1)$. Then,
 476 $\mathbb{E} [\|\mathbf{x}^t - \mathbf{x}^*\|^2] = O(1/t^\zeta)$, or equivalently, $\mathbb{E} [f(\mathbf{x}^t) - f^*] = O(1/t^\zeta)$. Here, $\zeta \in (0, 1)$
 477 is any positive number smaller than: $\min \left\{ 2\delta - 1, \frac{4a\mu(1-\kappa)\lambda(\kappa)(1-\delta)\mathcal{N}(1)}{L(aC'_2 + \|\mathbf{x}^0\| + \|\mathbf{x}^*\|) + B_0} \right\}$, where κ is
 478 an arbitrary constant in $(0, 1)$, and we recall quantities B_0 and $\lambda(\kappa)$ in Assumption 8;
 479 μ and L in Assumption 1; and C'_2 in Assumption 9. In alternative, let Assumptions 1,
 480 8, 9, 11, and 12 hold. Let $\alpha_t = \frac{a}{(t+1)^\delta}$, $\delta \in (0.5, 1]$, and assume that $\inf_{\mathbf{w} \neq 0} \frac{\|\Psi(\mathbf{w})\|}{\|\mathbf{w}\|} >$
 481 0 . Then, $\mathbb{E} [\|\mathbf{x}^t - \mathbf{x}^*\|^2] = O(1/t^\delta)$, or equivalently, $\mathbb{E} [f(\mathbf{x}^t) - f^*] = O(1/t^\delta)$. In
 482 particular, for $\delta = 1$ and a sufficiently large parameter a , we obtain the $O(1/t)$ MSE
 483 rate.*

484 **Example 4.1.** We illustrate the rate ζ in Theorem 4.2 for the gradient clipping
 485 nonlinearity with floor level $M > 0$. We consider an arbitrary joint pdf $p : \mathbb{R}^d \mapsto \mathbb{R}_+$
 486 that has “radial symmetry”, i.e., $p(\mathbf{u}) = q(\|\mathbf{u}\|)$, where $q : \mathbb{R}_+ \mapsto \mathbb{R}_+$ is a given

487 function. For example, we let:

488 (4.1)
$$p(\mathbf{u}) = q(\|\mathbf{u}\|), \quad q(\rho) = \frac{(\alpha - 2)(\alpha - 1)}{2\pi} \frac{1}{(1 + \rho)^\alpha}, \quad \rho \geq 0, \quad \alpha > 3.$$

489 It can be shown that $p(\mathbf{u})$ in (4.1) has finite moments of order r , $r < \alpha - 2$, and it
 490 has infinite moments for $r \geq \alpha - 2$. It holds that (see Appendix H for derivations)
 491 the rate ζ can be estimated as: $\min \{2\delta - 1, (1 - \delta) \frac{0.68\mu}{L}\}$. Hence, up to universal
 492 constants, the rate ζ is approximated as $\min \{2\delta - 1, (1 - \delta) \frac{\mu}{L}\}$. It is easy to see
 493 that the same rate estimate can be obtained for the normalized gradient nonlinearity,
 494 under the same gradient noise setting.

495 We compare the rate estimate here with the rate for component-wise nonlinearities
 496 (e.g., component-wise clipping in Example 3.3) that is, up to universal constants,
 497 of order $\min \{2\delta - 1, (1 - \delta) \frac{\mu}{\sqrt{d}L}\}$. We can see that, with the joint nonlinearity
 498 examples here, the rate is improved with respect to the component-wise nonlinearities
 499 by a factor \sqrt{d} . In other words, the rate estimate for the joint nonlinearities does
 500 not deteriorate with the dimension d increase. This may be intuitively explained
 501 by considering the sign component-wise nonlinearity and the normalized gradient.
 502 These two functions coincide for $d = 1$ (and this is reflected by the identical rate
 503 estimates we obtain here), but they become different for $d > 1$ (as also reflected by
 504 our obtained rate estimates). Intuitively, in the noiseless case, the normalized gradient
 505 preserves “more information” about the exact gradient (“true search direction”) than
 506 the component-wise sign function; hence, the difference in the estimated rates.

507 We now examine asymptotic normality for the joint nonlinearities case. We have
 508 the following theorem.

509 **THEOREM 4.3** (Asymptotic normality: Joint nonlinearity). *Consider algorithm*
 510 *(2.3) for solving optimization problem (2.1), and let Assumptions 1, 2, 8, 9, 10,*
 511 *and 13 hold. Further, let the step-size sequence $\{\alpha_t\}$ equal $\alpha_t = a/(t + 1)$, $a >$
 512 0 . Then: $\sqrt{t + 1}(\mathbf{x}^t - \mathbf{x}^*) \xrightarrow{d} \mathbb{N}(0, \mathcal{S})$. The asymptotic covariance \mathcal{S} is given by
 513 $\mathcal{S} = a^2 \int_0^\infty e^{v\mathbf{\Sigma}} \mathcal{S}_0 e^{v\mathbf{\Sigma}} dv$, where $\mathcal{S}_0 = \int \mathbf{u}\mathbf{u}^\top (\mathcal{N}(\|\mathbf{u}\|))^2 p(\mathbf{u}) d\mathbf{u}$; $\mathbf{\Sigma} = \frac{1}{2\mathbf{I} + a\mathbf{B}}$; $\mathbf{B} =$
 514 $-(\int \mathcal{N}(\|\mathbf{u}\|) p(\mathbf{u}) d\mathbf{u} + \int_{\mathbf{u} \neq 0} \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|} \mathcal{N}'(\|\mathbf{u}\|) p(\mathbf{u}) d\mathbf{u}) \nabla^2 f(\mathbf{x}^*)$, and constant $a > 0$ in the
 515 step-size sequence is taken large enough such that matrix $\mathbf{\Sigma}$ is stable. Moreover, the
 516 result continues to hold if Assumption 10 is replaced with Assumption 11, and As-
 517 sumption 12 is additionally imposed.*

518 Theorem 4.3 shows that asymptotic normality continues to hold for the joint nonlin-
 519 earity case as well, provided that $\mathcal{N}(a)$ is differentiable for any $a > 0$ and that \mathcal{N} is
 520 uniformly bounded from above.

521 **5. Intermediate results and proofs: Component-wise nonlinearities.**

522 This section provides proofs of Theorem 3.1, Theorem 3.2, and Theorem 3.3, ac-
 523 companied with the required intermediate results. Subsection 5.1 presents some use-
 524 ful intermediate results on stochastic approximation and deterministic time-varying
 525 sequences; Subsection 5.2 deals with the asymptotic analysis (Theorem 3.1 and The-
 526 orem 3.3); and Subsection 5.3 considers MSE analysis (Theorem 3.2).

527 **5.1. Stochastic approximation and time-varying sequences.** We present
 528 a useful result on single time scale stochastic approximation; see [26], Theorems 4.4.4
 529 and 6.6.1.

530 THEOREM 5.1. Let $\{\mathbf{x}^t \in \mathbb{R}^d\}$ be a random sequence that satisfies:

$$531 \quad (5.1) \quad \mathbf{x}^{t+1} = \mathbf{x}^t + \alpha_t [\mathbf{r}(\mathbf{x}^t) + \boldsymbol{\gamma}(t+1, \mathbf{x}^t, \omega)],$$

532 where, $\mathbf{r}(\cdot) : \mathbb{R}^d \mapsto \mathbb{R}^d$ is Borel measurable and $\{\boldsymbol{\gamma}(t, \mathbf{x}, \omega)\}_{t \geq 0, \mathbf{x} \in \mathbb{R}^d}$ is a family
533 of random vectors in \mathbb{R}^d , defined on a probability space $(\Omega, \mathcal{F}, \mathcal{P})$, and $\omega \in \Omega$ is a
534 canonical element. Let the following sets of assumptions hold:

535 **(B1)**: The function $\boldsymbol{\gamma}(t, \cdot, \cdot) : \mathbb{R}^d \times \Omega \mapsto \mathbb{R}^d$ is $\mathcal{B}^d \otimes \mathcal{F}$ measurable for every t ; \mathcal{B}^d is
536 the Borel algebra of \mathbb{R}^d .

537 **(B2)**: There exists a filtration $\{\mathcal{F}_t\}_{t \geq 0}$ of \mathcal{F} , such that, for each t , the family of ran-
538 dom vectors $\{\boldsymbol{\gamma}(t, \mathbf{x}, \omega)\}_{\mathbf{x} \in \mathbb{R}^d}$ is \mathcal{F}_t measurable, zero-mean and independent of \mathcal{F}_{t-1} .

539 **(B3)**: There exists a twice continuously differentiable function $V(\mathbf{x})$ with bounded
540 second order partial derivatives and a point $\mathbf{x}^* \in \mathbb{R}^d$ satisfying: $V(\mathbf{x}^*) = 0$, $V(\mathbf{x}) > 0$,
541 $\mathbf{x} \neq \mathbf{x}^*$, $\lim_{\|\mathbf{x}\| \rightarrow \infty} V(\mathbf{x}) = \infty$, $\sup_{\epsilon < \|\mathbf{x} - \mathbf{x}^*\| < \frac{1}{\epsilon}} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle < 0$, for any $\epsilon > 0$.

542 **(B4)**: There exist constants $k_1, k_2 > 0$, such that,

$$543 \quad \|\mathbf{r}(\mathbf{x})\|^2 + \mathbb{E} \left[\|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\|^2 \right] \leq k_1 (1 + V(\mathbf{x})) -$$

$$544 \quad - k_2 \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle.$$

546 **(B5)**: The weight sequence $\{\alpha_t\}$ satisfies $\alpha_t > 0$, $\sum_{t \geq 0} \alpha_t = \infty$, $\sum_{t \geq 0} \alpha_t^2 < \infty$.

547 **(C1)**: The function $\mathbf{r}(\mathbf{x})$ admits the representation

$$548 \quad (5.2) \quad \mathbf{r}(\mathbf{x}) = \mathbf{B}(\mathbf{x} - \mathbf{x}^*) + \boldsymbol{\delta}(\mathbf{x}),$$

549 where $\lim_{\mathbf{x} \rightarrow \mathbf{x}^*} \frac{\|\boldsymbol{\delta}(\mathbf{x})\|}{\|\mathbf{x} - \mathbf{x}^*\|} = 0$.

550 **(C2)**: The step-size sequence, $\{\alpha_t\}$ is of the form, $\alpha_t = \frac{a}{t+1}$, for any $t \geq 0$, where
551 $a > 0$ is a constant.

552 **(C3)**: Let \mathbf{I} be the $d \times d$ identity matrix and a, \mathbf{B} as in C2 and C1, respectively. Then,
553 the matrix $\Sigma = a\mathbf{B} + \frac{1}{2}\mathbf{I}$ is stable.

554 **(C4)**: The entries of the matrices, for any $t \geq 0, \mathbf{x} \in \mathbb{R}^d$, $\mathbf{A}(t, \mathbf{x}) = \mathbb{E}[\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)$
555 $\boldsymbol{\gamma}^\top(t+1, \mathbf{x}, \omega)]$ are finite, and the following limit exists: $\lim_{t \rightarrow \infty, \mathbf{x} \rightarrow \mathbf{x}^*} \mathbf{A}(t, \mathbf{x}) = \mathcal{S}_0$.

556 **(C5)**: There exists $\epsilon > 0$, such that

$$557 \quad (5.3) \quad \lim_{R \rightarrow \infty} \sup_{\|\mathbf{x} - \mathbf{x}^*\| < \epsilon} \sup_{t \geq 0} \int_{\|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\| > R} \|\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\|^2 dP = 0.$$

558 Then we have the following:

559 Let Assumptions **(B1)**-**(B5)** hold for $\{\mathbf{x}^t\}$ in (5.1). Then, starting from an
560 arbitrary initial state, the process $\{\mathbf{x}^t\}$ converges a.s. to \mathbf{x}^* .

561 The normalized process, $\{\sqrt{t}(\mathbf{x}^t - \mathbf{x}^*)\}$, is asymptotically normal if, besides As-
562 sumptions **(B1)**-**(B5)**, Assumptions **(C1)**-**(C5)** are also satisfied. In particular, as
563 $t \rightarrow \infty$, we have: $\sqrt{t}(\mathbf{x}^t - \mathbf{x}^*) \xrightarrow{d} \mathbb{N}(0, \mathcal{S})$. Also, the asymptotic covariance \mathcal{S} of the
564 multivariate Gaussian distribution $\mathbb{N}(0, \mathcal{S})$ is $\mathcal{S} = a^2 \int_0^\infty e^{v \Sigma} \mathcal{S}_0 e^{v \Sigma^\top} dv$.

565 *Proof.* For a proof see [26] (c.f. Theorems 4.4.4, 6.6.1). \square

566 We also make use of the following Theorem, proved in Appendix A; see also
567 Lemmas 4 and 5 in [21].

568 THEOREM 5.2. Let z^t be a nonnegative (deterministic) sequence satisfying:

$$569 \quad z^{t+1} \leq (1 - r_1^t) z^t + r_2^t,$$

570 for all $t \geq t'$, for some $t' > 0$, with some $z^{t'} \geq 0$. Here, $\{r_1^t\}$ and $\{r_2^t\}$ are deterministic
 571 sequences with $\frac{a_1}{(t+1)^{\delta_1}} \leq r_1^t \leq 1$ and $r_2^t \leq \frac{a_2}{(t+1)^{\delta_2}}$, with $a_1, a_2 > 0$, and $\delta_2 > \delta_1 > 0$.
 572 Then, the following holds: (1) If $\delta_1 < 1$, then $z^t = O(\frac{1}{t^{\delta_2 - \delta_1}})$; (2) If $\delta_1 = 1$, then
 573 $z^t = O(\frac{1}{t^{\delta_2 - 1}})$ provided that $a_1 > \delta_2 - \delta_1$; (3) if $\delta_1 = 1$ and $a_1 \leq \delta_2 - 1$, then
 574 $z^t = O(\frac{1}{t^\zeta})$, for any $\zeta < a_1$.

575 **5.2. Asymptotic analysis: Proofs of Theorem 3.1 and Theorem 3.3.** The
 576 next Lemma, due to [30], establishes structural properties of function ϕ in (3.1). The
 577 Lemma says that essentially, the convolution-like transformation of the nonlinearity
 578 preserves the structural properties of the nonlinearity. For a proof of the Lemma, see
 579 Appendix D.

580 **LEMMA 5.3.** [30] Consider function ϕ in (3.1), where function $\Psi : \mathbb{R} \mapsto \mathbb{R}$
 581 satisfies Assumption 5, and noise pdf $p : \mathbb{R} \mapsto \mathbb{R}_+$ satisfies Assumption 3. Then, the
 582 following holds.

- 583 1. ϕ is odd;
- 584 2. If in addition Assumption 7 holds, then $|\phi(a)| \leq K_2$, for any $a \in \mathbb{R}$, for some
 585 constant $K_2 > 0$;
- 586 3. If in addition Assumption 6 holds, then $|\phi(a)| \leq K_1(1 + |a|)$, for any $a \in \mathbb{R}$,
 587 for some constant $K_1 > 0$;
- 588 4. $\phi(a)$ is monotonically nondecreasing;
- 589 5. If in addition either Assumption 6 or Assumption 7 holds, then ϕ is differ-
 590 entiable at zero, with a strictly positive derivative at zero, equal to:

$$591 \quad (5.4) \quad \phi'(0) = \sum_{i=1}^s (\Psi(\nu_i + 0) - \Psi(\nu_i - 0)) p(\nu_i) + \sum_{i=0}^s \int_{\nu_i}^{\nu_{i+1}} \Psi'(\nu) p(\nu) d\nu,$$

592 where $\nu_i, i = 1, \dots, s$ are points of discontinuity of Ψ such that $\nu_0 = -\infty$ and
 593 $\nu_{s+1} = +\infty$.

594 **Remark.** In view of (5.4), we highlight the need that $p(u)$ is strictly positive in
 595 the vicinity of zero and that Ψ is either discontinuous at zero or strictly increasing
 596 in the vicinity of zero, in order for $\phi'(0)$ to be strictly positive. (see Assumptions 3
 597 and 5.) Consider the following counterexample: $\Psi(u) = \text{sign}(u)$, where p corresponds
 598 to the uniform distribution on the set $(-u_2, -u_1) \cup (u_1, u_2)$, for $0 < u_1 < u_2$. Note
 599 that p is zero in the vicinity of zero. Then, by (5.4), $\phi'(0) = 0$.

600 We proceed by setting up the proof of Theorem 3.1. The proof relies on conver-
 601 gence analysis of single-time scale stochastic approximation methods from [26]; more
 602 precisely, we utilize Theorem 5.1 in the Appendix; see also [20].

603 We first put algorithm (2.3) in the format that complies with Theorem 5.1.
 604 Namely, algorithm (2.3) can be written as:

$$605 \quad (5.5) \quad \mathbf{x}^{t+1} = \mathbf{x}^t + \alpha_t [\mathbf{r}(\mathbf{x}^t) + \gamma(t+1, \mathbf{x}^t, \omega)].$$

606 Here, ω denotes an element of the underlying probability space, and

$$607 \quad (5.6) \quad \mathbf{r}(\mathbf{x}) = -\phi(\nabla f(\mathbf{x})),$$

608 where, abusing notation, $\phi : \mathbb{R}^d \mapsto \mathbb{R}^d$ equals $(\phi(a_1, \dots, a_d)) = (\phi(a_1), \dots, \phi(a_d))^\top$.
 609 That is, we have that:

$$610 \quad (5.7) \quad \mathbf{r}(\mathbf{x}) = -(\phi([\nabla f(x)]_1), \dots, \phi([\nabla f(x)]_d))^\top, \quad \gamma(t+1, \mathbf{x}, \omega) = \phi(\nabla f(\mathbf{x})) - \Psi(\nabla f(\mathbf{x})) + \nu^t.$$

611 We provide an intuition behind the algorithmic format (5.5). Quantity $\mathbf{r}(\mathbf{x})$ is a
 612 deterministic, “useful”, progress direction with respect to the evolution of \mathbf{x}^t ; quantity
 613 $\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)$ is the stochastic component that plays a role of a noise in the system.

614 We adopt the following Lyapunov function: $V(x) = f(x) - f^*$, $V : \mathbb{R}^d \mapsto \mathbb{R}$,
 615 where $f^* = \inf_{x \in \mathbb{R}^d} f(x) = f(x^*)$. By Assumptions 1 and 2, V is twice continuously
 616 differentiable and has uniformly bounded second order partial derivatives, as required
 617 by Theorem 5.1. We are ready to prove Theorem 3.1.

618 *Proof* (Proof of Theorem 3.1). We now verify conditions B1-B5 from Theorem 5.1.
 619 Recall from Section 3 \mathcal{F}_t , the σ -algebra generated with random vectors $\boldsymbol{\nu}^s$, $s = 0, \dots, t$.
 620 Then, the family of random vectors $\{\boldsymbol{\gamma}(t+1, \mathbf{x}, \omega)\}_{\mathbf{x} \in \mathbb{R}^d}$ is \mathcal{F}_t -measurable, zero-mean
 621 and independent of \mathcal{F}_{t-1} . Also, clearly, function $\boldsymbol{\gamma}(t+1, \cdot, \cdot)$ is measurable, for all t .
 622 Thus, conditions B1 and B2 hold.

623 For B3, we need to prove that $\sup_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}^*\| \in (\epsilon, \frac{1}{\epsilon})} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle < 0$, for any $\epsilon > 0$,
 624 where $\nabla V(\mathbf{x}) = \nabla f(\mathbf{x})$. Let us fix an $\epsilon > 0$. Then, we have, for any $\mathbf{x} \in \mathbb{R}^d$:

$$\begin{aligned} 625 \quad \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle &= -\phi(\nabla f(\mathbf{x}))^\top (\nabla f(\mathbf{x})) \\ 626 \quad &= -\sum_{j=1}^d \phi([\nabla f(\mathbf{x})]_j) [\nabla f(\mathbf{x})]_j = -\sum_{j=1}^d |\phi([\nabla f(\mathbf{x})]_j)| |[\nabla f(\mathbf{x})]_j|, \end{aligned}$$

627 where the last equality holds because ϕ is an odd function. Consider arbitrary \mathbf{x} such
 628 that $\|\mathbf{x} - \mathbf{x}^*\| \geq \epsilon$. As $\|\nabla f(\mathbf{x})\|^2 \geq \mu^2 \|\mathbf{x} - \mathbf{x}^*\|^2$ (due to strong convexity of f), we
 629 have $\|\nabla f(\mathbf{x})\| \geq \mu\epsilon$, where we recall that μ is the strong convexity constant of f .
 630 Therefore, there exists an index $i \in \{1, \dots, d\}$ such that $|[\nabla f(\mathbf{x})]_i| \geq \frac{1}{d}\mu\epsilon =: \epsilon'$. Next,
 631 because $\phi'(0) > 0$, and ϕ is continuous at 0 and is non-decreasing (by Lemma 5.3),
 632 we have that $|\phi(b)| \geq \delta$ for some $\delta = \delta(\epsilon) > 0$, for all $b \in [\epsilon, 1/\epsilon]$. Finally, we
 633 have that: $\langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle \leq -\epsilon'\delta(\epsilon)$, for any \mathbf{x} such that $\|\mathbf{x} - \mathbf{x}^*\| \in [\epsilon, \frac{1}{\epsilon}]$, and
 634 therefore $\sup_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}^*\| \in (\epsilon, \frac{1}{\epsilon})} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle \leq \sup_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}^*\| \in [\epsilon, \frac{1}{\epsilon}]} \langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle \leq -\delta(\epsilon)\epsilon'$
 635 < 0 , hence verifying condition B3.

636 We next verify condition B4. Consider quantity $\mathbf{r}(\mathbf{x})$ in (5.6). By Lemma 5.3
 637 and the fact that f has Lipschitz gradient and is strongly convex (Assumption 1), it
 638 follows that: $\|\mathbf{r}(\mathbf{x})\|^2 \leq C_{r,1} + C_{r,2}V(\mathbf{x})$, for some positive constants $C_{r,1}$ and $C_{r,2}$.
 639 Also, since $\|\boldsymbol{\gamma}(\mathbf{x}, t+1, \omega)\|^2 \leq 2\|\phi(\nabla f(\mathbf{x}))\|^2 + 2\|\Psi(\nabla f(\mathbf{x}) + \boldsymbol{\nu}^t)\|^2$, and it holds that
 640 either 1) Ψ is bounded or 2) $|\Psi(a)| \leq C_2(1 + |a|)$ and ν_i^t has a finite variance, we
 641 have: $\mathbb{E}[\|\boldsymbol{\gamma}(\mathbf{x}, t+1, \omega)\|^2] \leq C_3 + C_4V(\mathbf{x})$, for some positive constants C_3, C_4 . Now,
 642 we finally have:

$$643 \quad \|\mathbf{r}(\mathbf{x})\|^2 + \mathbb{E}[\|\boldsymbol{\gamma}(\mathbf{x}, t+1, \omega)\|^2] \leq C_5 + C_6V(\mathbf{x}),$$

644 for some positive constants C_5, C_6 , and hence condition B4 holds for a constant $k_1 > 0$
 645 and $k_2 = 0$.⁶ Condition B5 holds by the choice of the step size sequence $\{\alpha_t\}$ in
 646 the Theorem statement. Summarizing, all conditions B1-B5 hold true, and hence
 647 $\mathbf{x}^t \rightarrow \mathbf{x}^*$, almost surely. \square

648 We continue by proving Theorem 3.3.

649 *Proof* (Proof of Theorem 3.3). We prove the Theorem by verifying conditions
 650 C1-C5 in Theorem 5.1. To verify condition C1, consider $\mathbf{r}(\mathbf{x})$ in (5.6) and note that,

⁶Note that the term $-\langle \mathbf{r}(\mathbf{x}), \nabla V(\mathbf{x}) \rangle$ in condition B4 of Theorem 5.1 equals $\langle \phi(\mathbf{x}), \nabla f(\mathbf{x}) \rangle$. This quantity is nonnegative, for any $\mathbf{x} \in \mathbb{R}^d$, and so k_2 can be taken to be any positive number. In other words, setting $k_2 = 0$ in B4 corresponds to a tighter inequality than the corresponding inequality for any $k_2 > 0$.

651 using the mean value theorem, it can be expressed as follows:

$$\begin{aligned}
652 \quad (5.8) \quad \mathbf{r}(\mathbf{x}) &= -\phi(\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*)) \\
&= -\phi \left(\underbrace{\left[\int_0^1 \nabla^2 f(\mathbf{x}^* + t(\mathbf{x} - \mathbf{x}^*)) dt \right]}_{\mathbf{H}} (\mathbf{x} - \mathbf{x}^*) \right) \\
&= -\phi(\mathbf{H}(\mathbf{x} - \mathbf{x}^*)) = -\phi'(0)\nabla^2 f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + \delta(\mathbf{x}),
\end{aligned}$$

653 where $\lim_{\mathbf{x} \rightarrow \mathbf{x}^*} \frac{\|\delta(\mathbf{x})\|}{\|\mathbf{x} - \mathbf{x}^*\|} = 0$. Hence, in the notation of [Theorem 5.1](#), we have that $\mathbf{B} =$
654 $-\phi'(0)\nabla^2 f(\mathbf{x}^*)$. Therefore, C1 holds. Also, C2 holds, by assumptions of [Theorem 3.3](#).
655 Now, we consider C3, which requires that the matrix $\Sigma = a\mathbf{B} + \frac{1}{2}\mathbf{I}$ is stable (all
656 its eigenvalues have negative real parts), where $\mathbf{B} = -\phi'(0)\nabla^2 f(\mathbf{x}^*)$. Note that
657 $\Sigma = \frac{1}{2}\mathbf{I} - a\phi'(0)\nabla^2 f(\mathbf{x}^*)$. Clearly, Σ is stable for large enough a , because the matrix
658 $\phi'(0)\nabla^2 f(\mathbf{x}^*)$ is positive definite. More precisely, Σ is stable for $a > 1/(2\mu\phi'(0))$.
659 Therefore, condition C3 holds, provided that $a > 1/(2\mu\phi'(0))$. We next consider
660 condition C4. In the notation of [Theorem 5.1](#), consider the following quantity:

$$\begin{aligned}
661 \quad \mathbf{A}(t, \mathbf{x}) &= \mathbb{E} [\gamma(t+1, \mathbf{x}, \omega)\gamma(t+1, \mathbf{x}, \omega)^\top] \\
662 &= \mathbb{E} \left[(\phi(\nabla f(\mathbf{x})) - \Psi(\nabla f(\mathbf{x}) + \boldsymbol{\nu}^t)) ((\phi(\nabla f(\mathbf{x})) - \Psi(\nabla f(\mathbf{x}) + \boldsymbol{\nu}^t))^\top) \right] \\
663 \quad (5.9) &= \mathbb{E} \left[(\phi(\nabla f(\mathbf{x})) - \Psi(\nabla f(\mathbf{x}) + \boldsymbol{\nu}^0)) ((\phi(\nabla f(\mathbf{x})) - \Psi(\nabla f(\mathbf{x}) + \boldsymbol{\nu}^0))^\top) \right] \\
664 \quad (5.10) &= \mathbb{E} [\gamma(1, \mathbf{x}, \omega)\gamma(1, \mathbf{x}, \omega)^\top].
\end{aligned}$$

665 Consider the set Ω^* of all outcomes $\omega \in \Omega$ such that Ψ is continuous at $\boldsymbol{\nu}^0(\omega)$. Clearly,
666 the set Ω^* has the probability one. For every $\omega \in \Omega^*$, we have $\Upsilon(\omega) := \lim_{t \rightarrow \infty, \mathbf{x} \rightarrow \mathbf{x}^*}$
667 $\gamma(1, \mathbf{x}, \omega)\gamma(1, \mathbf{x}, \omega)^\top = \Psi(\boldsymbol{\nu}^0)\Psi(\boldsymbol{\nu}^0)^\top$. Note that, for any $\epsilon > 0$, the random family
668 $\|\gamma(1, \mathbf{x}, \omega)\gamma(1, \mathbf{x}, \omega)^\top\|, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon$ is dominated by an integrable random variable.
669 (See ahead (5.12)–(5.13).) Therefore, by the dominated convergence theorem, and
670 the fact that the entries of $\boldsymbol{\nu}^0$ are mutually independent with pdf $p(u)$, we have that:

$$671 \quad (5.11) \quad \lim_{t \rightarrow \infty, \mathbf{x} \rightarrow \mathbf{x}^*} \mathbf{A}(t, \mathbf{x}) =: \mathcal{S}_0 = \mathbb{E} [\Psi(\boldsymbol{\nu}^0) \cdot \Psi(\boldsymbol{\nu}^0)^\top] = \sigma_\Psi^2 \cdot \mathbf{I},$$

672 where $\sigma_\Psi^2 = \int |\Psi(a)|^2 p(a) da$. Therefore, condition C4 holds. We finally verify condi-
673 tion C5. We follow the arguments analogous to those in [Theorem 10](#) in [\[20\]](#). Condi-
674 tion C5 means uniform integrability of the family $\{\|\gamma(t+1, \mathbf{x}, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon}$.
675 We have: $\|\gamma(t+1, \mathbf{x}, \omega)\|^2 \leq 2\|\phi(\nabla f(\mathbf{x}))\|^2 + 2\|\psi(\nabla f(\mathbf{x}) + \boldsymbol{\nu}^t)\|^2$. First, consider the
676 case when [Assumptions 6](#) and [4](#) hold. Then:

$$\begin{aligned}
677 \quad \|\gamma(t+1, \mathbf{x}, \omega)\|^2 &\leq C_7 + C_8\|\mathbf{x} - \mathbf{x}^*\|^2 + C_9\|\boldsymbol{\nu}^t\|^2 \\
678 \quad (5.12) &\leq C_7 + C_8\epsilon^2 + C_9\|\boldsymbol{\nu}^t\|^2,
\end{aligned}$$

679 for some positive constants C_7, C_8, C_9 . Consider next the family
680 $\{\tilde{\gamma}(t+1, \mathbf{x}, \omega)\}_{t=0,1,\dots, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon}$, with $\tilde{\gamma}(t+1, \mathbf{x}, \omega) = C_7 + C_8\epsilon^2 + C_9\|\boldsymbol{\nu}^t\|^2$. The family
681 $\{\tilde{\gamma}(t+1, \mathbf{x}, \omega)\}_{t=0,1,\dots, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon}$ is i.i.d. and hence it is uniformly integrable. The
682 family $\{\|\gamma(t+1, \mathbf{x}, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon}$ is dominated by
683 $\{\tilde{\gamma}(t+1, \mathbf{x}, \omega)\}_{t=0,1,\dots, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon}$ that is uniformly integrable, and hence
684 $\{\|\gamma(t+1, \mathbf{x}, \omega)\|^2\}_{t=0,1,\dots, \|\mathbf{x} - \mathbf{x}^*\| < \epsilon}$ is also uniformly integrable. Hence, C5 holds.

685 Now, let [Assumption 7](#) hold. Then:

$$686 \quad (5.13) \quad \|\widehat{\gamma}(t+1, \mathbf{x}, \omega)\|^2 \leq C_{10} + C_{11}\|\mathbf{x} - \mathbf{x}^*\|^2 \leq C_{10} + C_{11}\epsilon^2.$$

687 Consider the family $\{\widehat{\gamma}(t+1, \mathbf{x}, \omega)\}_{t=0,1,\dots,\|\mathbf{x}-\mathbf{x}^*\|<\epsilon}$, with $\widehat{\gamma}(t+1, \mathbf{x}, \omega) = C_{10} + C_{11}\epsilon^2$.
 688 The family $\{\widehat{\gamma}(t+1, \mathbf{x}, \omega)\}_{t=0,1,\dots,\|\mathbf{x}-\mathbf{x}^*\|<\epsilon}$ is uniformly integrable, and condition C5
 689 is verified analogously to the previous case. Summarizing, we have established that
 690 all conditions C1-C5 of [Theorem 5.1](#) hold true, thus the proof of [Theorem 3.3](#) \square .

691 **5.3. MSE analysis: Proof of [Theorem 3.2](#).** We start with the following
 692 Lemma that upper bounds $\|\nabla f(\mathbf{x}^t)\|$.

693 **LEMMA 5.4.** *Let Assumptions 1, 3, 5, and 7 hold. Further, let the step-size se-*
 694 *quence $\{\alpha_t\}$ be $\alpha_t = a/(t+1)^\delta$, $a > 0$, $\delta \in (0.5, 1)$. Then, for each $t = 1, 2, \dots$, we*
 695 *have, a.s.:*

$$696 \quad (5.14) \quad \|\nabla f(\mathbf{x}^t)\| \leq G_t := L \left(a C_2 \sqrt{d} \frac{t^{1-\delta}}{1-\delta} + \|\mathbf{x}^0 - \mathbf{x}^*\| \right).$$

697 *Proof.* Consider [\(2.3\)](#). Because the output of each component nonlinearity Ψ is
 698 bounded in the absolute value by C_2 ([Assumption 7](#)), we have, for each $t \geq 1$:

$$699 \quad \|\mathbf{x}^t - \mathbf{x}^*\| \leq \|\mathbf{x}^0 - \mathbf{x}^*\| + a \sqrt{d} C_2 \sum_{s=0}^{t-1} \frac{1}{(s+1)^\delta}$$

$$700 \quad (5.15) \quad \leq \|\mathbf{x}^0 - \mathbf{x}^*\| + a C_2 \sqrt{d} \left(\frac{t^{1-\delta}}{1-\delta} \right).$$

701 Next, because ∇f is L -Lipschitz, we have: $\|\nabla f(\mathbf{x}^t)\| \leq L \|\mathbf{x}^t - \mathbf{x}^*\|$. Applying this
 702 inequality to [\(5.15\)](#), the result follows. \square

703 We will also make use of the following Lemma.

704 **LEMMA 5.5.** *There exists a positive constant ξ such that, for any $t = 1, 2, \dots$, there*
 705 *holds, almost surely, for each $j = 1, \dots, d$, that: $|\phi([\nabla f(\mathbf{x}^t)]_j)| \geq |[\nabla f(\mathbf{x}^t)]_j| \frac{\phi'(0)\xi}{2G_t}$,*
 706 *where G_t is defined in [\(5.14\)](#).*

707 *Proof.* Consider function ϕ in [\(3.1\)](#). By [Lemma 5.3](#), we have that $\phi'(0) > 0$
 708 and ϕ is continuous at zero.⁷ Because ϕ is differentiable at zero, using first order
 709 Taylor series, there holds: $\phi(u) = \phi(0) + \phi'(0)u + h(u)u = \phi'(0)u + h(u)u$, $u \in \mathbb{R}$,
 710 where $h: \mathbb{R} \mapsto \mathbb{R}$ is a function such that $\lim_{u \rightarrow 0} h(u) = 0$. Due to the latter property
 711 of h , there exists a positive number ξ such that $|h(u)| \leq \frac{\phi'(0)}{2}$, for all $u \in [0, \xi]$.
 712 Using the latter bound, we obtain that $\phi(u) \geq \frac{1}{2}\phi'(0)u$, $u \in [0, \xi]$. Now, because ϕ is
 713 non-decreasing (by [Lemma 5.3](#)), it holds for any $a' > \xi$ that $\phi(a) \geq \frac{\phi'(0)\xi a}{2a'}$, for any
 714 $a \in [0, a']$. Consider now $\nabla f(\mathbf{x}^t)$. By [Lemma 5.4](#), we have that $\|\nabla f(\mathbf{x}^t)\| \leq G_t$, a.s.,
 715 and so, for any $j = 1, \dots, d$, $|[\nabla f(\mathbf{x}^t)]_j| \leq G_t$. Therefore, setting $a' = G_t$, the Lemma
 716 follows. \square

717 We are now ready to prove [Theorem 3.2](#).

718 *Proof (Proof of [Theorem 3.2](#)).* Consider algorithm [\(2.3\)](#) under Assumptions 1, 3,
 719 5, and 7. By the Lipschitz property of ∇f , we have, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, that:

$$720 \quad f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2,$$

⁷As ϕ is an odd function, for simplicity, in the proof we consider only nonnegative arguments of ϕ , while analogous analysis applies for negative arguments of ϕ .

721 and so, almost surely:

$$722 \quad (5.16) \quad \begin{aligned} f(\mathbf{x}^{t+1}) &\leq f(\mathbf{x}^t) + (\nabla f(\mathbf{x}^t))^\top (-\alpha_t \Psi(\nabla f(\mathbf{x}^t) + \boldsymbol{\nu}^t)) \\ &\quad + \frac{L}{2} \alpha_t^2 \|\Psi(\nabla f(\mathbf{x}^t) + \boldsymbol{\nu}^t)\|^2. \end{aligned}$$

723 Next, letting $\boldsymbol{\eta}^t = \Psi(\nabla f(\mathbf{x}^t) + \boldsymbol{\nu}^t) - \phi(\nabla f(\mathbf{x}^t))$, and using the fact that Ψ has bounded
724 outputs, we obtain:

$$725 \quad (5.17) \quad \begin{aligned} f(\mathbf{x}^{t+1}) &\leq f(\mathbf{x}^t) + (\nabla f(\mathbf{x}^t))^\top (-\alpha_t \phi(\nabla f(\mathbf{x}^t))) \\ &\quad + \frac{L}{2} \alpha_t^2 d^2 C_2^2 - \alpha_t (\nabla f(\mathbf{x}^t))^\top \boldsymbol{\eta}^t, \text{ a.s.} \end{aligned}$$

726 Recall filtration \mathcal{F}_t . Taking conditional expectation, and using that $\mathbb{E}[\boldsymbol{\eta}^t | \mathcal{F}_t] = 0$, we
727 get that, almost surely:

$$728 \quad (5.18) \quad \mathbb{E}[f(\mathbf{x}^{t+1}) | \mathcal{F}_t] \leq f(\mathbf{x}^t) - \alpha_t (\nabla f(\mathbf{x}^t))^\top \phi(\nabla f(\mathbf{x}^t)) + \frac{L}{2} \alpha_t^2 d^2 C_2^2.$$

729 Next, using [Lemma 5.5](#), and the fact that $\alpha_t = a/(t+1)^\delta$, we obtain that. a.s.:

$$730 \quad (5.19) \quad \mathbb{E}[f(\mathbf{x}^{t+1}) | \mathcal{F}_t] \leq f(\mathbf{x}^t) - \frac{c'}{(t+1)} \|\nabla f(\mathbf{x}^t)\|^2 + \frac{L a^2 d^2 C_2^2}{2 (t+1)^{2\delta}},$$

731 where $c' = \frac{a(1-\delta)\xi\phi'(0)}{2L(aC_2\sqrt{d} + \|\mathbf{x}^0 - \mathbf{x}^*\|)}$. Next, by strong convexity of f , we have that
732 $\|\nabla f(\mathbf{x}^t) - \nabla f(\mathbf{x}^*)\|^2 \geq 2\mu(f(\mathbf{x}^t) - f^*)$. Using the latter inequality, subtracting
733 f^* from both sides of the inequality, taking expectation, and applying [Theorem 5.2](#),
734 claims (2) and (3), we obtain the desired MSE rate result.

735 We next consider the case when [Assumption 7](#) is replaced with [Assumption 6](#) and
736 [Assumption 4](#) is additionally imposed. Following analogous arguments as in the first
737 part of the proof, it can be shown that, a.s.:

$$738 \quad (5.20) \quad \begin{aligned} \mathbb{E}[f(\mathbf{x}^{t+1}) | \mathcal{F}_t] &\leq f(\mathbf{x}^t) - \alpha_t \phi(\nabla f(\mathbf{x}^t))^\top \nabla f(\mathbf{x}^t) \\ &\quad + \frac{L}{2} \alpha_t^2 (C_{13} + C_{14} \mathbb{E}[\|\boldsymbol{\nu}^t\|^2 | \mathcal{F}_t]), \end{aligned}$$

739 for some positive constants C_{13}, C_{14} . Next, because $\inf_{a \neq 0} \frac{|\phi(a)|}{|a|} > 0$, we have that
740 $\phi(\nabla f(\mathbf{x}^t))^\top \nabla f(\mathbf{x}^t) \geq C_{15} \|\nabla f(\mathbf{x}^t)\|^2$, for some constant $C_{15} > 0$. Using the latter
741 bound in (5.20), subtracting f^* from both sides of the inequality, taking expectation,
742 and applying [Theorem 5.2](#), claim (1) and (2), the result follows. \square

743 **6. Intermediate results and proofs: Joint nonlinearities.** [Subsection 6.1](#)
744 provides the required intermediate results, while [Subsection 6.2](#) proves [Theorem 4.1](#).

745 **6.1. Intermediate results: Joint nonlinearities.** Recall function $\mathcal{N} : \mathbb{R}_+ \mapsto$
746 \mathbb{R}_+ in [Assumption 9](#). We first state and prove the following Lemma on the properties
747 of function \mathcal{N} .

748 **LEMMA 6.1.** *Under [Assumption 9](#), for any $\mathbf{x}, \mathbf{u} \in \mathbb{R}^d$, such that $\|\mathbf{u}\| > \|\mathbf{x}\|$, there*
749 *holds:*

$$750 \quad (6.1) \quad \begin{aligned} |\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) - \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)| &\leq \\ &\frac{\|\mathbf{x}\|}{\|\mathbf{u}\|} [\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)]. \end{aligned}$$

751 *Proof.* Fix a pair $\mathbf{x}, \mathbf{u} \in \mathbb{R}^d$, such that $\|\mathbf{u}\| > \|\mathbf{x}\|$, and assume without loss of
 752 generality that $\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) \geq \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)$. Then, (6.1) is equivalent to:

$$753 \quad (6.2) \quad (\|\mathbf{u}\| - \|\mathbf{x}\|)\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) \leq (\|\mathbf{u}\| + \|\mathbf{x}\|)\mathcal{N}(\|\mathbf{x} - \mathbf{u}\|).$$

754 Denote by $\rho = \|\mathbf{u}\|$. Notice that: $\rho - \|\mathbf{x}\| \leq \|\mathbf{x} + \mathbf{u}\| \leq \|\mathbf{x}\| + \|\mathbf{u}\| = \|\mathbf{x}\| + \rho$, and
 755 similarly, $\rho + \|\mathbf{x}\| \geq \|\mathbf{x} - \mathbf{u}\| \geq \rho - \|\mathbf{x}\|$. As \mathcal{N} is non-increasing, it follows that:
 756 $\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) \leq \mathcal{N}(\rho - \|\mathbf{x}\|)$, and $\mathcal{N}(\|\mathbf{x} - \mathbf{u}\|) \geq \mathcal{N}(\rho + \|\mathbf{x}\|)$. Now, we have:

$$757 \quad (6.3) \quad (\|\mathbf{u}\| - \|\mathbf{x}\|)\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) \leq (\rho - \|\mathbf{x}\|)\mathcal{N}(\rho - \|\mathbf{x}\|),$$

758 and similarly:

$$759 \quad (6.4) \quad (\|\mathbf{u}\| + \|\mathbf{x}\|)\mathcal{N}(\|\mathbf{x} - \mathbf{u}\|) \geq (\rho + \|\mathbf{x}\|)\mathcal{N}(\rho + \|\mathbf{x}\|).$$

760 By assumption, function $a \mapsto a\mathcal{N}(a)$, $a > 0$, is non-decreasing, and so $(\rho - \|\mathbf{x}\|)\mathcal{N}(\rho - \|\mathbf{x}\|)$
 761 $\leq (\rho + \|\mathbf{x}\|)\mathcal{N}(\|\mathbf{x}\| + \rho)$. Thus, combining (6.3) and (6.4), we have that (6.2)
 762 holds, which is in turn equivalent to the claim of the Lemma. \square

763 We now define map $\phi : \mathbb{R}^d \mapsto \mathbb{R}^d$, as follows. For a fixed (deterministic) point
 764 $\mathbf{w} \in \mathbb{R}^d$, we let:

$$765 \quad (6.5) \quad \phi(\mathbf{w}) = \int \Psi(\mathbf{w} + \mathbf{u})p(\mathbf{u})d\mathbf{u} = \mathbb{E}[\Psi(\mathbf{w} + \nu^0)],$$

766 where the expectation is taken with respect to the joint pdf of the gradient noise at
 767 any iteration t , e.g., $t = 0$. The map $\phi : \mathbb{R}^d \mapsto \mathbb{R}^d$ is, abusing notation, a counterpart
 768 of the component-wise map $\phi : \mathbb{R} \mapsto \mathbb{R}$ in (3.1). We have the following Lemma.

769 **LEMMA 6.2.** *Under Assumptions 8 and 9, the following holds:*

$$770 \quad (6.6) \quad \phi(\mathbf{x})^\top \mathbf{x} \geq 2(1 - \kappa)\|\mathbf{x}\|^2 \int_{\mathcal{J}(\mathbf{x})} \mathcal{N}(\|\mathbf{x}\| + \|\mathbf{u}\|)p(\mathbf{u})d\mathbf{u},$$

771 where $\mathcal{J}(\mathbf{x}) = \{\mathbf{u} : \frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\|\|\mathbf{x}\|} \in [0, \kappa]\}$, and κ is any constant in the interval $(0, 1)$.

772 *Proof.* Let us fix arbitrary $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{x} \neq 0$. As $\Psi(\mathbf{a}) = \mathbf{a}\mathcal{N}(\|\mathbf{a}\|)$, we have:

(6.7)

$$773 \quad \phi(\mathbf{x})^\top \mathbf{x} = \int_{\mathbf{u} \in \mathbb{R}^d} \underbrace{(\mathbf{x} + \mathbf{u})^\top \mathbf{x} \mathcal{N}(\|\mathbf{x} + \mathbf{u}\|)}_{:= \mathcal{M}(\mathbf{x}, \mathbf{u})} p(\mathbf{u})d\mathbf{u}$$

$$774 \quad (6.8) \quad = \int_{J_1(\mathbf{x}) = \{\mathbf{u} : \mathbf{u}^\top \mathbf{x} \geq 0\}} \mathcal{M}(\mathbf{x}, \mathbf{u})p(\mathbf{u})d\mathbf{u} + \int_{J_2(\mathbf{x}) = \{\mathbf{u} : \mathbf{u}^\top \mathbf{x} < 0\}} \mathcal{M}(\mathbf{x}, \mathbf{u})p(\mathbf{u})d\mathbf{u}.$$

775

776 Note also that there holds: $\mathcal{M}(\mathbf{x}, \mathbf{u}) = (\|\mathbf{x}\|^2 + \mathbf{u}^\top \mathbf{x})\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|)$; and $\mathcal{M}(\mathbf{x}, -\mathbf{u}) =$
 777 $(\|\mathbf{x}\|^2 - \mathbf{u}^\top \mathbf{x})\mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)$. Therefore, using the fact that $p(\mathbf{u}) = p(-\mathbf{u})$, for all
 778 $\mathbf{u} \in \mathbb{R}^d$, we obtain: $\phi(\mathbf{x})^\top \mathbf{x} = \int_{J_1(\mathbf{x})} \mathcal{M}_2(\mathbf{x}, \mathbf{u}) p(\mathbf{u})d\mathbf{u}$, where $\mathcal{M}_2(\mathbf{x}, \mathbf{u}) = [(\|\mathbf{x}\|^2 +$
 779 $\mathbf{u}^\top \mathbf{x})\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + (\|\mathbf{x}\|^2 - \mathbf{u}^\top \mathbf{x})\mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)]$. There holds:

$$780 \quad (6.9) \quad \mathcal{M}_2(\mathbf{x}, \mathbf{u}) \geq \|\mathbf{x}\|^2[\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)] -$$

$$- \|\mathbf{u}\|\|\mathbf{x}\|[\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) - \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)].$$

781 Since $\mathbf{u} \in J_1(\mathbf{x})$, there holds $\|\mathbf{x} + \mathbf{u}\| \geq \|\mathbf{x} - \mathbf{u}\|$. Now, using Lemma 6.1, we have:

$$782 \quad (6.10) \quad \begin{aligned} \mathcal{M}_2(\mathbf{x}, \mathbf{u}) &\geq \|\mathbf{x}\|^2 [\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)] - \\ &\quad \|\mathbf{u}\| \|\mathbf{x}\| \frac{\|\mathbf{x}\|}{\|\mathbf{u}\|} |\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)| = 0. \end{aligned}$$

783 Therefore, we have: $\mathcal{M}_2(\mathbf{x}, \mathbf{u}) \geq 0$, for any $\mathbf{u} \in J_1(\mathbf{x})$, $\|\mathbf{u}\| > \|\mathbf{x}\|$. Now, consider
 784 $\mathcal{J}(\mathbf{x}) = \{\mathbf{u} \in \mathbb{R}^d : \mathbf{u}^\top \mathbf{x} \geq 0, \frac{\mathbf{u}^\top \mathbf{x}}{\|\mathbf{u}\| \|\mathbf{x}\|} \in [0, \kappa]\}$, where $\kappa \in (0, 1)$. Let us consider
 785 $\mathbf{u} \in \mathcal{J}(\mathbf{x})$ such that $\|\mathbf{u}\| > \|\mathbf{x}\|$. Then, using Lemma 6.1, we get:

$$786 \quad (6.11) \quad \begin{aligned} \mathcal{M}_2(\mathbf{x}, \mathbf{u}) &\geq \|\mathbf{x}\|^2 [\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)] \\ &\quad - \|\mathbf{u}\| \|\mathbf{x}\| \kappa \underbrace{|\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) - \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)|}_{\geq 0} \\ &\geq (1 - \kappa) \|\mathbf{x}\|^2 (\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)). \end{aligned}$$

787 Now, consider $\mathbf{u} \in \mathcal{J}(\mathbf{x})$ such that $\|\mathbf{u}\| \leq \|\mathbf{x}\|$. Then, there holds:

$$788 \quad (6.12) \quad \begin{aligned} \mathcal{M}_2(\mathbf{x}, \mathbf{u}) &\geq \|\mathbf{x}\|^2 [\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)] - \\ &\quad \underbrace{\|\mathbf{u}\| \|\mathbf{x}\| \kappa}_{\leq \|\mathbf{x}\|} \underbrace{|\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)|}_{\geq 0} \\ &\geq (1 - \kappa) \|\mathbf{x}\|^2 (\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|)). \end{aligned}$$

789 where the last inequality holds due to the fact that $|a - b| \leq |a| + |b|$, for any $a, b \in \mathbb{R}$.
 790 Now, we have:

$$791 \quad (6.13) \quad \begin{aligned} \mathcal{M}_2(\mathbf{x}, \mathbf{u}) &\geq (1 - \kappa) \|\mathbf{x}\|^2 \underbrace{(\mathcal{N}(\|\mathbf{x} + \mathbf{u}\|) + \mathcal{N}(\|\mathbf{x} - \mathbf{u}\|))}_{\geq \mathcal{N}(\|\mathbf{x}\| + \|\mathbf{u}\|)} \\ &\geq 2(1 - \kappa) \|\mathbf{x}\|^2 \mathcal{N}(\|\mathbf{x}\| + \|\mathbf{u}\|), \text{ for any } \mathbf{u} \in \mathcal{J}(\mathbf{x}). \end{aligned}$$

792 From (6.13), we finally get:

$$793 \quad (6.14) \quad \begin{aligned} \phi(\mathbf{x})^\top \mathbf{x} &\geq \int_{\mathcal{J}(\mathbf{x})} 2(1 - \kappa) \|\mathbf{x}\|^2 \mathcal{N}(\|\mathbf{x}\| + \|\mathbf{u}\|) p(\mathbf{u}) d\mathbf{u} \\ &= 2(1 - \kappa) \|\mathbf{x}\|^2 \int_{\mathcal{J}(\mathbf{x})} \mathcal{N}(\|\mathbf{x}\| + \|\mathbf{u}\|) p(\mathbf{u}) d\mathbf{u}. \end{aligned} \quad \square$$

794 **LEMMA 6.3.** *Let Assumptions 1, 8, and Assumption 9 with condition 3. hold (the*
 795 *nonlinearity with bounded outputs case). Then, for each $t = 1, 2, \dots$, we have:*

$$796 \quad (6.15) \quad \|\nabla f(\mathbf{x}^t)\| \leq G'_t := L \left(a C'_2 \frac{t^{1-\delta}}{1-\delta} + \|\mathbf{x}^0 - \mathbf{x}^*\| \right).$$

797 *Proof.* The proof is analogous to the proof of Lemma 5.4. □

798 **6.2. Proofs of Theorems 4.1, 4.2, and 4.3: Joint nonlinearities.** We are
 799 now ready to prove the results for the joint nonlinearities case.

800 *Proof* (Proof of Theorem 4.1) We carry out the proof again by verifying con-
 801 ditions B1-B5 in Theorem 5.1. Algorithm (2.3) admits again the representation in
 802 Theorem 5.1 with

$$803 \quad (6.16) \quad \mathbf{r}(\mathbf{x}) = -\phi(\nabla f(\mathbf{x})), \quad \gamma(t + 1, \mathbf{x}, \omega) = \phi(\nabla f(\mathbf{x})) - \Psi(\nabla f(\mathbf{x}) + \nu^t).$$

804 Conditions B1 and B2 hold analogously to the proof of [Theorem 3.1](#). Condition
 805 B3 follows from [Lemma 6.2](#). Condition B4 holds analogously to the proof of [Theo-](#)
 806 [rem 3.1](#). Finally, condition B5 follows from the definition of the step-size sequence in
 807 [Theorem 4.1](#). Thus, the result. \square We next prove [Theorem 4.3](#). *Proof* (Proof of [Theo-](#)
 808 [rem 4.3](#)) We carry out the proof again by verifying conditions C1–C5 in [Equation \(8.2\)](#).
 809 The conditions C2–C5 are verified analogously as in the proof of [Theorem 3.3](#). For
 810 condition C1, first fix an arbitrary $\mathbf{u} \neq 0$, and consider points \mathbf{x} in the vicinity of \mathbf{x}^* .
 811 Then, using the differentiability of $\mathcal{N}(a)$ for $a \neq 0$ and the differentiability of ∇f , it
 812 can be shown that:

$$813 \quad \Psi(\mathbf{u} + \nabla f(\mathbf{x})) = \mathbf{u}\mathcal{N}(\|\mathbf{u}\|) + \mathcal{N}(\|\mathbf{u}\|)\nabla^2 f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \\ 814 \quad + \mathcal{N}'(\|\mathbf{u}\|) \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|} \nabla^2 f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) + o(\|\mathbf{x} - \mathbf{x}^*\|).$$

815 We next integrate the above equality with respect to the joint pdf $p(\mathbf{u})$. For the first
 816 term above, note that $\int \mathcal{N}(\|\mathbf{u}\|)\mathbf{u}p(\mathbf{u})d\mathbf{u} = 0$, because $p(\mathbf{u}) = p(-\mathbf{u})$, for all \mathbf{u} . The
 817 second term is integrable as $\sup_{a>0} \mathcal{N}(a) < \infty$ ([Assumption 13](#)). The third term is
 818 integrable as function $a \mapsto a\mathcal{N}(a)$ is by assumptions non-decreasing; then, by taking
 819 its derivative, it follows that $|\mathcal{N}'(a)| \leq \mathcal{N}(a)/a$, $a > 0$, and so $\|\mathbf{u}\mathbf{u}^\top \mathcal{N}'(\|\mathbf{u}\|)\|/\|\mathbf{u}\|$
 820 $\leq \mathcal{N}(\|\mathbf{u}\|)$. Now, using the definition of $\mathbf{r}(\mathbf{x})$, it follows that $\mathbf{r}(\mathbf{x})$ admits the repre-
 821 sentation [\(5.2\)](#), with:

$$822 \quad \mathbf{B} = - \left(\int \mathcal{N}(\|\mathbf{u}\|)p(\mathbf{u})d\mathbf{u} + \int_{\mathbf{u} \neq 0} \frac{\mathbf{u}\mathbf{u}^\top}{\|\mathbf{u}\|} \mathcal{N}'(\|\mathbf{u}\|)p(\mathbf{u})d\mathbf{u} \right) \nabla^2 f(\mathbf{x}^*).$$

823 The conditions C1–C5 hold; thus, the result. \square

824 We are now ready to prove [Theorem 4.2](#).

825 *Proof* (Proof of [Theorem 4.2](#)) We first consider the case when [Assumptions 1, 8,](#)
 826 [9, and 10](#) hold. Analogously to the proof of [3.2](#), it can be shown that, a.s.:

$$827 \quad (6.17) \quad \mathbb{E}[f(\mathbf{x}^{t+1}) | \mathcal{F}_t] \leq f(\mathbf{x}^t) - \alpha_t \phi(\nabla f(\mathbf{x}^t))^\top \nabla f(\mathbf{x}^t) + \alpha_t^2 C_{17},$$

828 for some positive constant C_{17} . By [Lemma 6.2](#), there holds, for $\mathbf{a} := \nabla f(\mathbf{x}^t)$, a.s.:

$$829 \quad (6.18) \quad (\phi(\mathbf{a}))^\top \mathbf{a} \geq 2(1 - \kappa)\|\mathbf{a}\|^2 \int_{\mathcal{J}} \mathcal{N}(\|\mathbf{a}\| + \|\mathbf{u}\|)p(\mathbf{u})d\mathbf{u},$$

830 where we recall $\mathcal{J} = \{\mathbf{u} : \frac{\mathbf{u}^\top \mathbf{a}}{\|\mathbf{u}\|\|\mathbf{a}\|} \in [0, \kappa]\}$, and $\kappa \in (0, 1)$ is a constant. Note that, as

831 $a \mapsto a\mathcal{N}(a)$ is non-decreasing, \mathcal{N} satisfies: $\mathcal{N}(b) \geq \min\left(\frac{\mathcal{N}(1)}{b}, \mathcal{N}(1)\right)$ for any $b > 0$.

832 Consider constant B_0 in condition 2. of [Assumption 8](#). Then, for all \mathbf{u} such that

833 $\|\mathbf{u}\| \leq B_0$, there holds $\mathcal{N}(\|\mathbf{a}\| + \|\mathbf{u}\|) \geq \min\left\{\frac{\mathcal{N}(1)}{\|\mathbf{a}\| + B_0}, \mathcal{N}(1)\right\}$. We now have, a.s.:

$$834 \quad (6.19) \quad \|\nabla f(\mathbf{x}^t)\|^2 \int_{\mathcal{J}} \mathcal{N}(\|\nabla f(\mathbf{x}^t)\| + \|\mathbf{u}\|)p(\mathbf{u})d\mathbf{u}$$

$$835 \quad (6.20) \quad \geq \|\nabla f(\mathbf{x}^t)\|^2 \int_{\mathcal{J}_4} \min\left\{\frac{\mathcal{N}(1)}{B_0 + \|\nabla f(\mathbf{x}^t)\|}, \mathcal{N}(1)\right\}p(\mathbf{u})d\mathbf{u}$$

$$836 \quad (6.21) \quad \geq \|\nabla f(\mathbf{x}^t)\|^2 \frac{\mathcal{N}(1)}{B_0 + G'_t} \int_{\mathcal{J}_4} p(\mathbf{u})d\mathbf{u}.$$

837 Here, $J_4 = \{u \in \mathbb{R}^d : \frac{\mathbf{u}^\top \nabla f(\mathbf{x}^t)}{\|\mathbf{u}\| \|\nabla f(\mathbf{x}^t)\|} \in [0, \kappa], \|\mathbf{u}\| \leq B_0\}$. In (6.20), we used the fact that
 838 $\mathcal{N}(a)$ is non-negative for any $a \geq 0$, and in (6.21), we used Lemma 6.3.
 839 Therefore, we have that, almost surely, for sufficiently large t :

$$840 \quad \|\nabla f(\mathbf{x}^t)\|^2 \int_J \mathcal{N}(\|\nabla f(\mathbf{x}^t)\| + \|\mathbf{u}\|) p(\mathbf{u}) d\mathbf{u} \geq C_{18} \frac{\|\nabla f(\mathbf{x}^t)\|^2}{G'_t + B_0},$$

841 for some positive constant C_{18} .
 842 Combining the last bound with Lemmas 6.2 and 6.3, in view of condition 2. in
 843 Assumption 8, we obtain that, for sufficiently large t , a.s.:

$$844 \quad (6.22) \quad (\phi(\nabla f(\mathbf{x}^t)))^\top \nabla f(\mathbf{x}^t) \geq C_{19} \frac{\|\nabla f(\mathbf{x}^t)\|^2}{B_0 + G'_t},$$

845 where the positive constant C_{19} can be taken as $C_{19} = 2(1 - \kappa)\lambda(\kappa)\mathcal{N}(1)$. Applying
 846 the bound (6.22) to (6.17) we obtain an equivalent to (5.19). Therein, c' in (5.19) is
 847 replaced with a positive constant c'' that can be taken as $c'' = \frac{4a(1-\kappa)\lambda(\kappa)(1-\delta)\mathcal{N}(1)}{L(aC'_2 + \|\mathbf{x}^0 - \mathbf{x}^*\|) + B_0}$.
 848 We now proceed analogously to the proof of Theorem 3.2, by applying claims (2) and
 849 (3) of Theorem 5.2. The desired MSE result now follows, with the rate ζ being any
 850 positive number less than

$$851 \quad (6.23) \quad \min \left\{ 2\delta - 1, \frac{4a\mu(1-\kappa)\lambda(\kappa)(1-\delta)\mathcal{N}(1)}{L(aC'_2 + \|\mathbf{x}^0\| + \|\mathbf{x}^*\|) + B_0} \right\}.$$

852 We now consider the case when Assumptions 1, 8, 9, 11, and 12 hold. We have,
 853 by assumption, that $\inf_{\mathbf{x} \neq 0} \frac{\|\Psi(\mathbf{x})\|}{\|\mathbf{x}\|} > 0$. This is equivalent to saying that \mathcal{N} is lower-
 854 bounded by a positive constant, i.e., $\mathcal{N}(a) \geq C_{20}$, for each a , for some constant
 855 $C_{20} > 0$. Then, it follows that, a.s.:

$$856 \quad (6.24) \quad (\phi(\nabla f(\mathbf{x}^t)))^\top \nabla f(\mathbf{x}^t) \geq C_{21} \|\nabla f(\mathbf{x}^t)\|^2,$$

857 for some positive constant C_{21} . The proof then proceeds analogously to the proof of
 858 Theorem 3.2 by applying the appropriate variant of Theorem 5.2. \square

859 **7. Experiments.** In order to benchmark the proposed nonlinear SGD frame-
 860 work, we consider `Heart`, `Diabetes` and `Australian` datasets from the LibSVM li-
 861 brary [9]. We consider the logistic regression loss function for binary classification,
 862 see, e.g., [15], where function f in (2.1) is the empirical loss, i.e., the sum of the logistic
 863 losses across all data points in a given dataset.

864 As it has been studied in [15] (see Figure 2 in [15]), we have, near the solution \mathbf{x}^* ,
 865 the following behavior with respect to gradient noise. (See also [15] for details how
 866 the gradient noise is evaluated in Figure 2 therein.) With the `HEART` dataset, tails of
 867 stochastic gradients are not heavy. On the other hand, for `DIABETES` and `AUSTRALIAN`
 868 datasets, the gradient noise has outliers and exhibits a heavy-tail behavior.

869 We consider three different nonlinearities to demonstrate the effectiveness of our
 870 nonlinear framework, namely, `tanh` (hyperbolic tangent), `sign` and a bi-level cus-
 871 tomization of `sign` with $\Psi(x) = -1, -0.5, 0.5, 1$, for $x \in (-\infty, -0.5], (-0.5, 0],$
 872 $(0, 0.5], (0.5, \infty]$, respectively (`nonlinear-quantizer` in figures). Note that the `tanh`
 873 function may be considered a smooth approximation of `sign`. We benchmark the
 874 above methods against the linear SGD, clipped-SGD and SSTM along with a clipped
 875 version of SSTM from [15]. For each of the methods, we use batch sizes of 50, 100
 876 and 20 for the `Australian`, `Diabetes` and `Heart` datasets, respectively. We also

877 consider clipped-SGD with periodically decreasing clipping level (**d-clipped-SGD** in
 878 Figures) as a baseline as introduced in [15]. This method starts with some initial clip-
 879 ping level and after every l epochs the clipping level is multiplied by some constant
 880 $c \in (0, 1)$. The step sizes α_t (learning rates) for each method from our framework
 881 were tuned after an experimentation. The learning rates for the baselines, i.e., SGD,
 882 clipped-SGD, SSTM and clipped-SSTM are also tuned and are selected to be as in
 883 [15]. In more detail, the learning rates for the proposed methods are of the form
 884 $a/(b(t+1)+L)$, where we recall that t is the iteration counter, L is the smoothness
 885 constant of ∇f , and parameters a, b are tuned via grid search. The value of a is
 886 chosen to be 1.0, 1.5 and 5.0, respectively, for **Heart**, **Diabetes** and **Australian** and
 887 for all the three non-linearities. The value of b is chosen to be 0.001, 7.0 and 7.0
 888 respectively for **Australian**, **Heart** and **Diabetes** datasets for the **sign** nonlinearity.
 889 The value of b is chosen to be 0.0001, 2.0 and 3.0×10^{-6} respectively for **Australian**,
 890 **Heart** and **Diabetes** datasets for the **tanh** nonlinearity. The value of b is chosen to
 891 be 0.001, 5.0 and 5.0 respectively for **Australian**, **Heart** and **Diabetes** datasets for
 892 the **nonlinear-quantizer** nonlinearity.

893 We first note that (see Figure 3) **d-clipped-SGD** stabilizes the trajectory as com-
 894 pared to the linear SGD, even if the initial clipping level was high. At the same time,
 895 clipped-SGD with large clipping levels performs similarly as SGD. It is noteworthy,
 896 that SGD has the least oscillations for **Australian** and **Diabetes** datasets, despite
 897 the fact that these datasets have heavier or similar tails. This can be attributed to
 898 the fact that SGD does not get close to the solution in terms of functional value.
 899 SSTM in particular shows large oscillations, which can be attributed to it being a
 900 version of accelerated/momentum-based methods and its usage of small batch sizes.
 901 **Clipped-SSTM** on the other hand suffers less from oscillations and has a comparable
 902 convergence rate as SSTM. In comparison, all the three nonlinear schemes that have
 903 been proposed in this paper, have very little oscillations. While the **tanh** algorithm
 904 is outperformed by the algorithms with other nonlinearities from our framework, its per-
 905 formance is at par with the other baselines from [15]. In particular, the **sign** algorithm
 906 compares favorably to other baselines in terms of convergence for **Australian** and
 907 **Heart** datasets. The **nonlinear-quantizer** algorithm outperforms other baselines for
 908 the **Diabetes** dataset. The good behavior of **tanh** and **sign** on the heavy-tail data
 909 sets, specially relative to the linear SGD, also viewing **tanh** as a smooth approxima-
 910 tion of **sign**, might also be related with the insights from Example 3.4. In summary,
 911 the three simple example nonlinearities from the proposed framework are comparable
 912 or favorable over the considered state-of-the-art benchmarks on the studied datasets.
 913

914 **8. Conclusion.** We proposed a general framework for nonlinear stochastic gra-
 915 dient descent (SGD) under heavy-tail gradient noise. Unlike existing studies of SGD
 916 under heavy-tail noise that focus on specific nonlinear functions (e.g., adaptive clip-
 917 ping), our framework includes a broad class of component-wise (e.g., sign gradient)
 918 and joint (e.g., gradient clipping) nonlinearities. We establish for the considered meth-
 919 ods almost sure convergence, MSE convergence rate, and also asymptotic covariance
 920 for component-wise nonlinearities. We carry out numerical experiments on several real
 921 datasets that exhibit heavy tail gradient noise effects. The experiments show that,
 922 while our framework is more general than existing studies of SGD under heavy-tail
 923 noise, several easy-to-implement nonlinearities from our framework are competitive
 924 with state-of-the-art alternatives.

925

REFERENCES

- 926 [1] D. ALISTARH, D. GRUBIC, J. LI, R. TOMIOKA, AND M. VOJNOVIC, *QSGD: Communication-*
 927 *efficient sgd via gradient quantization and encoding*, in Advances in Neural Information
 928 Processing Systems, 2017, pp. 1709–1720.
- 929 [2] L. BALLE, F. PEDREGOSA, AND N. L. ROUX, *The geometry of sign gradient descent*, arXiv
 930 preprint arXiv:2002.08056, (2020).
- 931 [3] H. BERCOVICI AND V. PATA, *Stable laws and domains of attraction in free probability theory*,
 932 *Ann. of Math.*, 149 (1999), pp. 1023–1060.
- 933 [4] J. BERNSTEIN, Y.-X. WANG, K. AZIZZADENESHELI, AND A. ANANDKUMAR, *signsgd: Compressed*
 934 *optimisation for non-convex problems*, in International Conference on Machine Learning,
 935 PMLR, 2018, pp. 560–569.
- 936 [5] L. BOTTOU, *Large-scale machine learning with stochastic gradient descent*, in Proceedings of
 937 COMPSTAT’2010, Springer, 2010, pp. 177–186.
- 938 [6] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine*
 939 *learning*, *Siam Review*, 60 (2018), pp. 223–311.
- 940 [7] R. H. BYRD, G. M. CHIN, W. NEVEITT, AND J. NOCEDAL, *On the use of stochastic hessian*
 941 *information in optimization methods for machine learning*, *SIAM Journal on Optimization*,
 942 21 (2011), pp. 977–995.
- 943 [8] V. CEVHER, S. BECKER, AND M. SCHMIDT, *Convex optimization for big data: Scalable, ran-*
 944 *domized, and parallel algorithms for big data analytics*, *IEEE Signal Processing Magazine*,
 945 31 (2014), pp. 32–43.
- 946 [9] C.-C. CHANG AND C.-J. LIN, *Libsvm: a library for support vector machines*, *ACM transactions*
 947 *on intelligent systems and technology (TIST)*, 2 (2011), pp. 1–27.
- 948 [10] S. DASARATHAN, C. TEPEDELENLIOĞLU, M. K. BANAVAR, AND A. SPANIAS, *Robust consensus*
 949 *in the presence of impulsive channel noise*, *IEEE Transactions on Signal Processing*, 63
 950 (2015), pp. 2118–2129.
- 951 [11] D. DAVIS, D. DRUSVYATSKIY, L. XIAO, AND J. ZHANG, *From low probability to high confidence*
 952 *in stochastic convex optimization.*, *J. Mach. Learn. Res.*, 22 (2021), pp. 49–1.
- 953 [12] S. GHADIMI AND G. LAN, *Optimal stochastic approximation algorithms for strongly convex*
 954 *stochastic composite optimization i: A generic algorithmic framework*, *SIAM Journal on*
 955 *Optimization*, 22 (2012), pp. 1469–1492.
- 956 [13] S. GHADIMI AND G. LAN, *Optimal stochastic approximation algorithms for strongly convex*
 957 *stochastic composite optimization I: A generic algorithmic framework*, *SIAM J. Optim.*,
 958 22 (2012), pp. 1469–1492.
- 959 [14] S. GHADIMI AND G. LAN, *Optimal stochastic approximation algorithms for strongly convex*
 960 *stochastic composite optimization, II: Shrinking procedures and optimal algorithms*, *SIAM*
 961 *J. Optim.*, 23 (2013), pp. 2061–2089.
- 962 [15] E. GORBUNOV, M. DANILOVA, AND A. GASNIKOV, *Stochastic optimization with heavy-tailed*
 963 *noise via accelerated gradient clipping*, arXiv preprint arXiv:2005.10785, (2020).
- 964 [16] E. GORBUNOV, F. HANZELY, AND P. RICHTÁRIK, *A unified theory of sgd: Variance reduction,*
 965 *sampling, quantization and coordinate descent*, in International Conference on Artificial
 966 Intelligence and Statistics, PMLR, 2020, pp. 680–690.
- 967 [17] M. GURBUZBALABAN, U. SIMSEKLI, AND L. ZHU, *The heavy-tail phenomenon in sgd*, in Inter-
 968 national Conference on Machine Learning, PMLR, 2021, pp. 3964–3975.
- 969 [18] S. HORVÁTH, D. KOVALEV, K. MISHCHENKO, S. STICH, AND P. RICHTÁRIK, *Stochastic*
 970 *distributed learning with gradient quantization and variance reduction*, arXiv preprint
 971 arXiv:1904.05115, (2019).
- 972 [19] A. JUDITSKY, A. NAZIN, A. NEMIROVSKY, AND A. TSYBAKOV, *Algorithms of robust stochastic*
 973 *optimization based on mirror descent method*, arXiv:1907.02707, (2019).
- 974 [20] S. KAR, J. M. MOURA, AND K. RAMANAN, *Distributed parameter estimation in sensor net-*
 975 *works: Nonlinear observation models and imperfect communication*, *IEEE Transactions*
 976 *on Information Theory*, 58 (2012), pp. 3575–3605.
- 977 [21] S. KAR AND J. M. F. MOURA, *Convergence rate analysis of distributed gossip (linear parameter)*
 978 *estimation: Fundamental limits and tradeoffs*, *IEEE Jour. Sel. Top. Sig. Proc.*, 5 (2011),
 979 pp. 674–690.
- 980 [22] U. A. KHAN, S. KAR, AND J. M. MOURA, *Distributed average consensus: Beyond the realm of*
 981 *linearity*, in 2009 Conference Record of the Forty-Third Asilomar Conference on Signals,
 982 Systems and Computers, IEEE, 2009, pp. 1337–1342.
- 983 [23] L. LEI AND M. I. JORDAN, *On the adaptivity of stochastic gradient-based optimization*, *SIAM*
 984 *Journal on Optimization*, 30 (2020), pp. 1473–1500.
- 985 [24] H. MANIA, X. PAN, D. PAPALIOPOULOS, B. RECHT, K. RAMCHANDRAN, AND M. I. JORDAN,

- 986 *Perturbed iterate analysis for asynchronous stochastic optimization*, SIAM Journal on Op-
 987 timization, 27 (2017), pp. 2202–2229.
- 988 [25] A. NEMIROVSKI, A. JUDITSKY, G. LAN, AND A. SHAPIRO, *Robust stochastic approximation*
 989 *approach to stochastic programming*, SIAM Journal on optimization, 19 (2009), pp. 1574–
 990 1609.
- 991 [26] M. B. NEVELSON AND R. Z. KHASHINSKIĬ, *Stochastic approximation and recursive estimation*,
 992 vol. 47, American Mathematical Soc., 1976.
- 993 [27] F. NIU, B. RECHT, C. RÉ, AND S. J. WRIGHT, *Hogwild!: A lock-free approach to parallelizing*
 994 *stochastic gradient descent*, arXiv preprint arXiv:1106.5730, (2011).
- 995 [28] R. PASCANU, T. MIKOLOV, AND Y. BENGIO, *On the difficulty of training recurrent neural*
 996 *networks*, in International Conference on Machine Learning, PMLR, 2013, pp. 1310–1318.
- 997 [29] V. PICHAPATI, A. T. SURESH, F. X. YU, S. J. REDDI, AND S. KUMAR, *Adaclip: Adaptive*
 998 *clipping for private sgd*, arXiv preprint arXiv:1908.07643, (2019).
- 999 [30] B. T. POLYAK AND Y. Z. TSYPKIN, *Adaptive estimation algorithms: convergence, optimality,*
 1000 *stability*, Avtomatika i Telemekhanika, (1979), pp. 71–84.
- 1001 [31] A. SHAPIRO, D. DENTCHEVA, AND A. RUSZCZYNSKI, *Lectures on stochastic programming: mod-*
 1002 *eling and theory*, SIAM, 2021.
- 1003 [32] U. SIMSEKLI, M. GÜRBÜZBALABAN, T. H. NGUYEN, G. RICHARD, AND L. SAGUN, *On the*
 1004 *heavy-tailed theory of stochastic gradient descent for deep neural networks*, arXiv preprint
 1005 arXiv:1912.00018, (2019).
- 1006 [33] S. S. STANKOVIĆ, M. BEKO, AND M. S. STANKOVIĆ, *A robust consensus seeking algorithm*, in
 1007 IEEE EUROCON 2019-18th International Conference on Smart Technologies, IEEE, 2019,
 1008 pp. 1–6.
- 1009 [34] S. SUNDARAM AND B. GHARESIFARD, *Consensus-based distributed optimization with malicious*
 1010 *nodes*, in 2015 53rd Annual Allerton Conference on Communication, Control, and Com-
 1011 puting (Allerton), IEEE, 2015, pp. 244–249.
- 1012 [35] F. YOUSEFIAN, A. NEDIĆ, AND U. V. SHANBHAG, *On stochastic gradient and subgradient meth-*
 1013 *ods with adaptive steplength sequences*, Automatica, 48 (2012), pp. 56–67.
- 1014 [36] J. ZHANG, T. HE, S. SRA, AND A. JADBABAIE, *Why gradient clipping accelerates training: A*
 1015 *theoretical justification for adaptivity*, arXiv preprint arXiv:1905.11881, (2019).
- 1016 [37] J. ZHANG, S. P. KARIMIREDDY, A. VEIT, S. KIM, S. J. REDDI, S. KUMAR, AND S. SRA, *Why*
 1017 *are adaptive methods good for attention models?*, arXiv preprint arXiv:1912.03194, (2019).

1018 Appendix.

1019 **A. Proof of Theorem 5.2.** We first state and prove the following Lemma.

1020 LEMMA 8.1. *Consider (deterministic) sequence*

$$1021 \quad v^{t+1} = \left(1 - \frac{a_3}{(t+1)^\delta}\right) v^t + \frac{a_4}{(t+1)^\delta}, \quad t \geq t_0,$$

1022 *with $a_3, a_4 > 0$ and $0 < \delta \leq 1$, $t_0 > 0$, and $v^{t_0} \geq 0$. Further, assume that t_0 is such*
 1023 *that $\frac{a_3}{(t+1)^\delta} \leq 1$, for all $t \geq t_0$. Then, $\lim_{t \rightarrow \infty} v^t = \frac{a_4}{a_3}$.*

1024 *Proof.* Let $e^t = v^t - \frac{a_4}{a_3}$. It is easy to verify that:

$$1025 \quad e^{t+1} = \left(1 - \frac{a_3}{(t+1)^\delta}\right) e^t, \quad t \geq t_0.$$

1026 Then, for all $t \geq t_0$, there holds:

$$1027 \quad (8.1) \quad |e^{t+1}| = \left(1 - \frac{a_3}{(t+1)^\delta}\right) |e^t| \leq \exp\left(-a_3 \sum_{s=t_0}^t \frac{1}{(s+1)^\delta}\right) |e^{t_0}|$$

1028 where in (8.1) we used the inequality $1 + a \leq \exp(a)$, $a > 0$. Letting $t \rightarrow \infty$ and the
 1029 fact that $\delta \leq 1$ so that the sequence $\frac{1}{(s+1)^\delta}$, $s \geq t_0$, is non-summable, we obtain that
 1030 $e^t \rightarrow 0$, which in turn implies the claim of the Lemma. \square

1031 We now continue with proving [Theorem 5.2](#). First, let us prove claim (1). Note
1032 that:

$$1033 \quad (8.2) \quad z^{t+1} \leq \left(1 - \frac{a_1}{(t+1)^{\delta_1}}\right) z^t + \frac{a_2}{(t+1)^{\delta_2}}, \quad t \geq t'.$$

1034 Multiplying the above inequality with $(t+1)^{\delta_2-\delta_1}$, defining $\widehat{z}^t = t^{\delta_2-\delta_1} z^t$, we get:

$$1035 \quad \widehat{z}^{t+1} \leq \left(1 - \frac{a_1}{(t+1)^{\delta_1}}\right) (1+1/t)^{\delta_2-\delta_1} \widehat{z}^t + \frac{a_2}{(t+1)^{\delta_1}}.$$

1036 Next, using, e.g., a Taylor expansion of function $a \mapsto (1+a)^{\delta_2-\delta_1}$, it can be shown
1037 that $(1+1/t)^{\delta_2-\delta_1} \leq 1 + \frac{2(\delta_2-\delta_1)}{t}$, for any $t \geq t_\delta$, for appropriately chosen $t_\delta > 0$.
1038 Therefore,

$$1039 \quad \left(1 - \frac{a_1}{(t+1)^{\delta_1}}\right) (1+1/t)^{\delta_2-\delta_1}$$

$$1040 \quad \leq 1 - \frac{a_1}{(t+1)^{\delta_1}} + \frac{2(\delta_2-\delta_1)}{t} - \frac{2a_1(\delta_2-\delta_1)}{t(t+1)^{\delta_1}} \leq 1 - \frac{a_1}{2(t+1)^{\delta_1}},$$

1041 for any $t \geq t_1$, for appropriately taken $t_1 > 0$. Using the latter bound, we obtain:
1042 $\widehat{z}^{t+1} \leq \left(1 - \frac{a_1}{2(t+1)^{\delta_1}}\right) \widehat{z}^t + \frac{a_2}{(t+1)^{\delta_1}}$, $t \geq t_1$. Now, applying [Lemma 8.1](#), we obtain that
1043 $\widehat{z}^t = O(1)$, and therefore $z^t = O(1/t^{\delta_2-\delta_1})$. This proves claim (1) in [Theorem 5.2](#).

1044 We now prove claim (2). Multiplying (8.2) by $(t+1)^{\delta_2-1}$, and defining $\widehat{z}^t =$
1045 $t^{\delta_2-1} z^t$, we obtain:

$$1046 \quad \widehat{z}^{t+1} \leq \left(1 - \frac{a_1}{(t+1)}\right) (1+1/t)^{\delta_2-1} \widehat{z}^t + \frac{a_2}{t+1}$$

$$1047 \quad (8.3) \quad \leq \left(1 - \frac{a_1 - (\delta_2 - 1)}{t} + \frac{C_{22}}{t^2}\right) \widehat{z}^t + \frac{a_2}{t+1}$$

$$1048 \quad (8.4) \quad \leq \left(1 - \frac{a_1 - (\delta_2 - 1)}{2(t+1)}\right) \widehat{z}^t + \frac{a_2}{t+1}, \quad t \geq t_2,$$

1049 for appropriately chosen $t_2 > 0$ and $C_{22} > 0$. In (8.3), we used the fact that $(1 +$
1050 $1/t)^{\delta_2-1} \leq 1 + \frac{\delta_2-1}{t} + \frac{C_{23}}{t^2}$, for all $t \geq 1$ and some $C_{23} > 0$ (the inequality can
1051 be obtained, e.g., via a Taylor approximation). The claim (2) of [Theorem 5.2](#) now
1052 follows by applying [Lemma 8.1](#) to (8.4).

1053 We now prove claim (3). Let $a_1 < \delta_2 - 1$, and fix an arbitrary positive number ζ ,
1054 $\zeta < a_1$. Then, we have, for $\widehat{z}^t = t^\zeta z^t$:

$$1055 \quad \widehat{z}^{t+1} \leq \left(1 - \frac{a_1}{(t+1)}\right) (1+1/t)^\zeta \widehat{z}^t + \frac{a_2}{(t+1)^{\delta_2-\zeta}}$$

$$1056 \quad \leq \left(1 - \frac{a_1 - \zeta}{t} + \frac{C_{24}}{t^2}\right) \widehat{z}^t + \frac{a_2}{(t+1)^{\delta_2-\zeta}}$$

$$1057 \quad \leq \left(1 - \frac{a_1 - \zeta}{2(t+1)}\right) \widehat{z}^t + \frac{a_2}{t+1}, \quad t \geq t_3,$$

1058 for appropriately chosen $t_3 > 0$ and $C_{24} > 0$. In the last inequality, we used the fact
1059 that $\zeta < a_1 \leq \delta_2 - 1$, and so $\delta_2 - \zeta > 1$. Finally, applying [Lemma 8.1](#), claim (3)
1060 follows. \square

1061 **B. A demonstration that the linear SGD's iterate sequence has infinite**
 1062 **variance.** We provide here a simple demonstration that the linear SGD's iterate
 1063 sequence has infinite variance under the setting of [Assumption 1](#), [Assumption 3](#), and
 1064 [Assumption 5](#), condition 3., holds.

1065 More precisely, assume that the gradient noise ν^t has infinite variance. Consider
 1066 algorithm (2.3) for solving problem (1) with $f : \mathbb{R} \mapsto \mathbb{R}$, $f(x) = \frac{x^2}{2}$, with Ψ being
 1067 the identity function. Further, consider arbitrary sequence of positive step-sizes $\{\alpha_t\}$.
 1068 Then, we have:

$$1069 \quad (8.5) \quad x^{t+1} = (1 - \alpha_t) x^t - \alpha_t \nu^t, \quad t = 0, 1, \dots,$$

1070 with arbitrary deterministic initialization $x^0 \in \mathbb{R}$. Then, squaring (8.5), using the
 1071 independence of x^t and ν^t , and the fact that ν^t has zero mean, we get: $\mathbb{E}[(x^{t+1})^2]$
 1072 $= (1 - \alpha_t)^2 \mathbb{E}[(x^t)^2] + \alpha_t^2 \mathbb{E}[(\nu^t)^2] \geq \alpha_t^2 \mathbb{E}[(\nu^t)^2]$, $t = 0, 1, \dots$ Taking expectation and
 1073 using the fact that $\mathbb{E}[(\nu^t)^2] = +\infty$, we see that $\mathbb{E}[(x^t)^2] = +\infty$, for any $t \geq 1$.

1074 **C. Extension of [Theorem 3.2](#) for gradient noise vector with mutually**
 1075 **dependent entries.** We show that [Theorem 3.2](#) continues to hold when we have an
 1076 i.i.d. zero mean noise vector sequence $\{\nu^t\}$ with a joint pdf $p : \mathbb{R}^d \mapsto \mathbb{R}$. In more
 1077 detail, we provide an extension of [Lemma 6.2](#) but for component-wise nonlinearities.

1078 Namely, as in [Lemma 6.2](#), consider, for a fixed $\mathbf{y} \neq 0$:

$$1079 \quad (8.6) \quad \int \psi(\mathbf{y} + \mathbf{u})^\top \mathbf{y} p(\mathbf{u}) \, d\mathbf{u}.$$

1080 As, for $\mathbf{a} \in \mathbb{R}^d$, we have $\Psi(\mathbf{a}) = (\Psi(a_1), \dots, \Psi(a_d))^\top$ (component-wise nonlinearity),
 1081 we have:

$$1082 \quad \int \psi(\mathbf{y} + \mathbf{u})^\top \mathbf{y} p(\mathbf{u}) \, d\mathbf{u} = \int \left(\sum_{i=1}^d \psi(y_i + u_i) y_i \right) p(\mathbf{u}) \, d\mathbf{u}$$

$$1083 \quad = \sum_{i=1}^d \int (\psi(y_i + u_i) y_i) p(\mathbf{u}) \, d\mathbf{u} = \sum_{i=1}^d \int (\psi(y_i + u_i) y_i) p_i(u_i) \, du_i,$$

1084 where $p_i(u_i)$ is the marginal pdf of the i -th component of ν^t . It is easy to show,
 1085 as $p(\mathbf{u}) = p(-\mathbf{u})$, $\mathbf{u} \in \mathbb{R}^d$, that, for any $i = 1, \dots, d$, we have $p_i(u) = p_i(-u)$, $u \in$
 1086 \mathbb{R} . Define $\phi_i(a) = \int \Psi(a + u) p_i(u) \, du$. Note that $\phi_i(a)$ now obeys [Lemma 5.3](#). In
 1087 particular, ϕ_i is also odd, and hence:

$$1088 \quad \int \psi(\mathbf{y} + \mathbf{u})^\top \mathbf{y} p(\mathbf{u}) \, d\mathbf{u} = \sum_{i=1}^d \int (\psi(y_i + u_i) y_i) p_i(u_i) \, du_i$$

$$1089 \quad = \sum_{i=1}^d \phi_i(y_i) y_i = \sum_{i=1}^d |\phi_i(y_i)| |y_i|.$$

1090 The last inequality holds because, for any $i = 1, \dots, d$, quantities $\phi_i(y_i)$ and y_i have
 1091 equal sign. The proof now proceeds analogously to that of [Theorem 3.2](#).

1092 **D. Proof of [Lemma 5.3](#).** The proof can be found in [30]; we include similar
 1093 arguments for completeness. For claim 1., note that

$$1094 \quad \phi(a) = \int_{-\infty}^{+\infty} \Psi(a + u) p(u) \, du = - \int_{-\infty}^{+\infty} \Psi(-a - u) p(u) \, du$$

$$1095 \quad = - \int_{-\infty}^{+\infty} \Psi(-a + w) p(w) \, dw = -\phi(-a),$$

1096 for any $a \in \mathbb{R}$, where we use the fact that Ψ is odd. For claim 2., note that $|\phi(a)| \leq$
 1097 $\int_{-\infty}^{+\infty} |\Psi(a+u)|p(u)du \leq C_1 \int_{-\infty}^{+\infty} p(u)du = C_1$, where we used [Assumption 7](#). Proof
 1098 of claim 3. is similar to that of claim 2. For claim 4., note that $\phi(a) = \int_0^{+\infty}$
 1099 $(\Psi(u+a) - \Psi(u-a))p(u)du$, and so, for $a' > a$, we have

$$1100 \quad \phi(a') - \phi(a) = \int_0^{+\infty} [(\Psi(u+a') - \Psi(u+a)) +$$

$$1101 \quad + (\Psi(u-a) - \Psi(u-a'))]p(u)du \geq 0,$$

1102 because Ψ is non-decreasing. Finally, for claim 5., to show that $\phi'(0)$ is given by [\(5.4\)](#),
 1103 see the proof of Lemma 6 in [\[30\]](#). To verify that $\phi'(0)$ is strictly positive, consider first
 1104 the case that Ψ has a discontinuity at zero. Then, because $p(0) > 0$ by [Assumption 3](#),
 1105 it follows from [\(5.4\)](#) that $\phi'(0) \geq (\Psi(0+) - \Psi(0-))p(0) > 0$. Otherwise, if Ψ is
 1106 continuous at zero, we have: $\phi'(0) \geq \int_{-c}^c \Psi'(u)p(u)du > 0$, where $c > 0$ is taken
 1107 such that $\Psi(u)$ is continuous and strictly increasing and $p(u)$ is strictly positive for
 1108 $|u| < c$.⁸ Such c exists in view of Assumptions [3](#) and [5](#).

1109 **E. Derivations for Example 3.3.** We calculate the rate ζ in [Theorem 3.2](#) for
 1110 the component-wise clipping nonlinearity with saturation value m , $m > 1$. Here, it
 1111 can be shown, by doing direct calculations, that

$$1112 \quad (8.7) \quad \phi(w) = 2w \int_0^{m-w} p(u)du + \int_{m-w}^{m+w} (m+w-u)p(u)du, \quad w \in [0, m].$$

1113 Furthermore, it can be shown that (see Appendix F): $\phi'(0) = 2 \int_0^m p(u)du$. Noting
 1114 that the second integral in [\(8.7\)](#) is nonnegative, and using the form $p(u)$ in [\(3.2\)](#), we
 1115 obtain:

$$1116 \quad (8.8) \quad \phi(w) \geq 2w \int_0^{m-w} p(u)du = w \left(1 - \frac{1}{(m-w+1)^{\alpha-1}} \right), \quad w \in [0, m].$$

1117 Also, we have: $\phi'(0) = 1 - \frac{1}{(m+1)^\alpha}$. From the latter equation and [\(8.8\)](#), we estimate
 1118 that ξ can be taken as: $\xi = m+1 - \left(\frac{2}{1+(m+1)^{-\alpha}} \right)^{1/(\alpha-1)} \geq m-1$, for any $\alpha > 2$, for
 1119 any $m > 1$. Hence, we can also take $\xi = m-1$. Substituting the obtained estimates
 1120 for $\phi'(0)$ and ξ into the rate ζ , we obtain the rate estimate in [\(3.3\)](#).

1121 **F. Derivation of $\phi'(0)$ for Example 3.5.** Consider the coordinate-wise clipping
 1122 nonlinearity Ψ with floor level $m > 0$. The function Ψ here is piece-wise differentiable,
 1123 with the derivative $\Psi'(a) = 1$, for $a \in (-m, m)$, and $\Psi'(a) = 0$, for $|a| > m$. We now
 1124 apply claim 5. in [Lemma 5.3](#) and use formula [\(5.4\)](#) for evaluating $\phi'(0)$. As the
 1125 coordinate-wise clipping function does not have discontinuity points, [\(5.4\)](#) simplifies
 1126 to the following:

$$1127 \quad \phi'(0) = \int_{u \in \mathbb{R}, u \neq -m, u \neq m} \Psi'(u)p(u) du = \int_{-m}^{+m} p(u)du = 2 \int_0^m p(u)du,$$

1128 where the last equality uses symmetry of function $p(u)$.

⁸If there are some (at most countably many) points inside interval $(-c, c)$ where Ψ is continuous but not differentiable, these points are excluded from the integration set in $\int_{-c}^c \Psi'(u)p(u)du$ without change in the integration result.

1129 **G. Derivations for Example 3.6.** We provide here details for the derivations
1130 in Example 3.6. We first calculate σ_{Ψ}^2 ; we have:

$$1131 \quad \sigma_{\Psi}^2 = \int_{-\infty}^{\infty} |\Psi(u)|^2 p(u) du = \int_{-\infty}^{\infty} p(u) du = 1.$$

1132 Next, by direct integration, we have for $\alpha > 3$:

$$1133 \quad \sigma_{\nu}^2 = 2 \int_0^{\infty} p(u) u^2 du$$

$$1134 \quad = -(\alpha - 1) \frac{[(\alpha - 1)u((\alpha - 2)u + 2)] + 2}{(\alpha - 3)(\alpha - 2)(\alpha - 1)(1 + u)^{\alpha - 2}} \Big|_0^{\infty} = \frac{2}{(\alpha - 3)(\alpha - 2)}.$$

1135 On the other hand, for $\alpha \in (2, 3]$, we clearly have $\sigma_{\nu}^2 = +\infty$. Finally, using claim 5.
1136 in [Lemma 5.3](#), and using the fact that $\Psi'(u) = 0$, for all $u \neq 0$, we obtain:

$$1137 \quad \phi'(0) = p(0) (\Psi(0+) - \Psi(0-)) = 2p(0) = \alpha - 1.$$

1138 **H. Derivations for Example 4.1.** We consider the (joint) gradient clipping
1139 nonlinearity Ψ with the clipping level $M > 0$, and we consider $p(\mathbf{u})$ in (4.1).

1140 Consider rate ζ in [Theorem 4.2](#) that, for a sufficiently large a , can be approximated
1141 as:

$$1142 \quad (8.9) \quad \min \left\{ 2\delta - 1, (1 - \delta) \frac{4\mu(1 - \kappa)\lambda(\kappa)\mathcal{N}(1)}{L C_2'} \right\}.$$

1143 Here, κ is an arbitrary scalar in $(0, 1)$, and, for the gradient clipping, we have that
1144 $\mathcal{N}(1) = C_2' = M$. Note that, regarding [Assumption 8](#), quantity B_0 can be taken here
1145 to be an arbitrary positive number. Moreover, for $p(\mathbf{u})$ in (4.1), due to the radial
1146 symmetry, we have that

$$1147 \quad \lambda(\kappa) = \lambda(\kappa, B_0) = \frac{1}{\pi} \arccos(1 - \kappa) \mathcal{P}(B_0), \quad \kappa \in (0, 1),$$

1148 where $\mathcal{P}(B_0) = \int_{\mathbf{u}: \|\mathbf{u}\| \leq B_0} p(\mathbf{u}) d\mathbf{u} = 1 - \frac{1 + (\alpha - 1)B_0}{(1 + B_0)^{\alpha - 1}}$. We next maximize (8.10), i.e., we
1149 maximize $(1 - \kappa)\lambda(\kappa, B_0)$ with respect to $\kappa \in (0, 1)$, to get the largest (tightest) esti-
1150 mate of ζ . It is easy to see that $\max_{\kappa \in (0, 1)} (1 - \kappa)\lambda(\kappa, B_0) > 0.17 \mathcal{P}(B_0)$. Substituting
1151 all the above developments into (8.10), we obtain:

$$1152 \quad (8.10) \quad \zeta \approx \min \left\{ 2\delta - 1, (1 - \delta) \frac{0.68\mu \mathcal{P}(B_0)}{L} \right\}$$

$$1153 \quad = \min \left\{ 2\delta - 1, (1 - \delta) \frac{0.68\mu}{L} \left(1 - \frac{1 + (\alpha - 1)B_0}{(1 + B_0)^{\alpha - 1}} \right) \right\}$$

1154 As B_0 can be arbitrary positive number, letting $B_0 \rightarrow +\infty$, we obtain the following
1155 rate estimate: $\min \left\{ 2\delta - 1, (1 - \delta) \frac{0.68\mu}{L} \right\}$. It is easy to see that the same rate esti-
1156 mate can be obtained for the normalized gradient nonlinearity. The only difference
1157 in the rate derivation is that therein $\mathcal{N}(1) = C_2' = 1$.